

Week2

2017年9月11日 16:37

I Introduction

1) What's machine learning

① two problems

- { classification / pattern recognition
- detection / hypothesis testing
- regression / prediction / estimation / filtering

classification

example, image recognition

input, image

output, "what is the image"

mathematically,

$$\begin{array}{ccc} x & \xrightarrow{\quad} & \boxed{f(\cdot)} \xrightarrow{\quad} c \\ \text{real vector} & & \\ x \in \mathbb{R}^n & & c \in \{c_1, c_2, \dots, c_r\} = C \\ f: \mathbb{R}^n \rightarrow C & & \text{↑ people playing frisbee} \end{array}$$

estimation / regression .

Example, price prediction / forecasting

input, price at $t, t+1, t+2, \dots, t+n$ (real vector)

Output, price at $t+k$ ($k > 0$) (real number)

Mathematically,

$$\begin{array}{ccc} x & \xrightarrow{\quad} & \boxed{f(\cdot)} \xrightarrow{\quad} y \\ \text{real vector} & & \\ x \in \mathbb{R}^n & & y \in \mathbb{R}^m \\ f: \mathbb{R}^n \rightarrow \mathbb{R}^m & & \end{array}$$

② Machine learning

To get $f(\cdot)$ \rightarrow Design $f(\cdot)$

Δ requires understanding
(theory)

\rightarrow By machine learning

Δ requires "data" training data.

D. training data
M.L. D. $\rightarrow f(\cdot)$

2) Mathematical foundation

- Probability theory & statistics
- Optimization

2. Review . Prob. RV & statistics

1) Probability

① Prob. model , Prob Axioms , Prob. space

Sample space
universe

Probability model

σ -Algebra. \curvearrowleft

Experiment
Outcome
Event
Probability

Ex: Roll a dice once

$\{1, 2, 3 \dots 6\} = \Sigma$
"get an Even"
 $\{2, 4, 6\}$

$P["\text{get an event } \#"] = P[\{2, 4, 6\}] = \frac{1}{2}$

Event,

An event A is subset of S Events, $A, B \dots \rightarrow f$ set of event.If A is an event, A^c must be an eventIf A & B are event, So are $A \cap B, A \cup B$.If A_1, A_2, \dots, A_n are events, \rightarrow So is $\bigcup_{i=1}^n A_i$

Sigma.

 σ -Algebra. A^c — A complement

Probability: Measure, normalized.

A measure, function, $m(A)$

input: a set

output: non-negative number

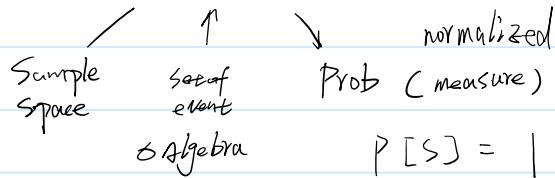
additivity: if $A \cap B = \emptyset$, $m(A \cup B) = m(A) + m(B)$

Normalized

$$0 \leq P[A] \leq 1, \quad P[S] = 1$$

Summary, Probability model Probability space

$$(S, f, P)$$

Ex₂

Experiment: roll a dice once

$$S = \{1, 2, 3, \dots, 6\}$$

 $f =$ all subsets of S .

$$P, P[\{\cdot\}] = \dots /$$

Sum of normally distributed random variables

$$X \sim N(\mu_X, \sigma_X^2)$$

$$Y \sim N(\mu_Y, \sigma_Y^2)$$

$$Z = X + Y,$$

then

$$Z \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2).$$

 $S =$ same

$$f_1 = \{ \{2, 4, 6\} \}$$

expand f_1 such that f_1 is σ -Algebra

↳ minimally

specify a prob for this

quiz

2017年9月11日 16:37

3 偏导函数的定义

如果 $z = f(x, y)$ 在区域 D 内任一点 (x, y) 处对 x 的偏导数都存在，那么这个偏导数就是 x 、 y 的函数，称为 $z = f(x, y)$ 对 x 的偏导（函）数，

记作 $\frac{\partial z}{\partial x}$, $\frac{\partial f}{\partial x}$ 或 $f'_x(x, y)$

$z = f(x, y)$ 对 y 的偏导数

记作 $\frac{\partial z}{\partial y}$, $\frac{\partial f}{\partial y}$, 或 $f'_y(x, y)$

<https://wenku.baidu.com/view/e4e6659e51e79b8968022.html#rc-view>

注 1 如何计算偏导数

在计算 $\frac{\partial z}{\partial x}$ 将 y 看成常量，仅对 x 求导

在计算 $\frac{\partial z}{\partial y}$ 将 x 看成常量，仅对 y 求导

注 2 一点的偏导数与偏导函数的关系

$$\left. \frac{\partial z}{\partial x} \right|_{\substack{x=x_0 \\ y=y_0}} = \frac{\partial z}{\partial x} \text{ 在 } (x_0, y_0) \text{ 点的函数值}$$

A set of subsets of Ω , \mathcal{F} , is called a **σ -field (σ -algebra)** if it satisfies the following 3 properties:

1. $\Omega \in \mathcal{F}$
2. \mathcal{F} is closed under complementation: If $A \in \mathcal{F}$, then so is its complement, $\Omega \setminus A = A^c \in \mathcal{F}$.
3. \mathcal{F} is closed under countable unions: If $A, B \in \mathcal{F}$, then so is $A \cup B \in \mathcal{F}$.

This means that $A \cap B = (A^c \cup B^c)^c$ from de Morgan's law = closed under intersections.

or $(A \cup B)^c = (A^c \cap B^c)$

The power set is the largest σ -algebra of Ω .

Trivial set $\leq \sigma$ -algebra $\leq P(\Omega)$

An algebra is a collection of subsets which is closed under some countable set operations i.e. if a set X can be obtained from subsets in \mathcal{F} through complementation, union or intersection then X is also a member of \mathcal{F} .

Sigma Field / sigma algebra

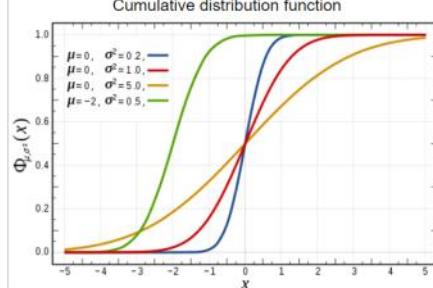
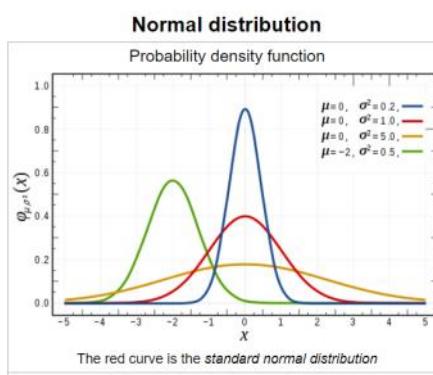
根据组合数和二项式定理

子集个数: $C_0 + C_1 + C_2 + \dots + C_n = (1+1)^n = 2^n$

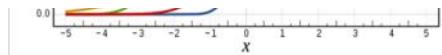
<https://www.math.lsu.edu/~sengupta/7312s02/sigmaalg.html>

normal (or Gaussian) distribution

Notation	$\mathcal{N}(\mu, \sigma^2)$
Parameters	$\mu \in \mathbb{R}$ — mean (location) $\sigma^2 > 0$ — variance (squared scale)
Support	$x \in \mathbb{R}$
PDF	$\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
CDF	$\frac{1}{2} \left[1 + \operatorname{erf} \left(\frac{x-\mu}{\sigma\sqrt{2}} \right) \right]$
Quantile	$\mu + \sigma\sqrt{2} \operatorname{erf}^{-1}(2F - 1)$
Mean	μ
Median	μ
Mode	μ
Variance	σ^2
Skewness	0
Ex. kurtosis	0
Entropy	$\frac{1}{2} \log(2\pi e \sigma^2)$



Skewness	0
Ex. kurtosis	0
Entropy	$\frac{1}{2} \log(2\pi e \sigma^2)$
MGF	$\exp\{\mu t + \frac{1}{2}\sigma^2 t^2\}$
CF	$\exp\{i\mu t - \frac{1}{2}\sigma^2 t^2\}$
Fisher information	$\begin{pmatrix} 1/\sigma^2 & 0 \\ 0 & 1/(2\sigma^4) \end{pmatrix}$



<https://www.zhihu.com/question/26055805>

Week3

2017年9月18日 16:37

1) Joint & Conditional probs.

$$\text{Joint: } P[A \cap B]$$

↑
joint

$$\text{Conditional: } P[A|B] = \frac{P[A \cap B]}{P[B]}$$

$$\Rightarrow \begin{cases} \text{Bayesian influence} \\ P[A \cap B] = P[A|B] \cdot P[B] \end{cases}$$

? ← independence
 $P[A|B] = P[A]$
 $P[A \cap B] = P[A] \cdot P[B]$

- Reverse influence.
 $P[B|A] = \frac{P[A|B] \cdot P[B]}{P[A]}$

cause
obs.
cause
obs.

- Total prob formula.

$$P[A] = P[A \cap B] + P[A \cap B^c]$$

$$= P[A|B] \cdot P[B] + P[A|B^c] \cdot P[B^c]$$

2) Random Variables.

A random variable:

① basic idea: convert (S, f, P) into (\mathbb{R}, B, P)

Example: flip a coin once.

$$(S, f, P) \longrightarrow (\mathbb{R}, B, P)$$

$\#2$

$$S = \{H, T\}$$

$$f = \{f(H), f(T), f(H, T), \emptyset\} \quad X \text{ is a measurable function. } \text{Corresponds to } f.$$

$$P[\{H\}] = \frac{1}{2}$$

$P(\text{Event}) = \int \text{pdf}$

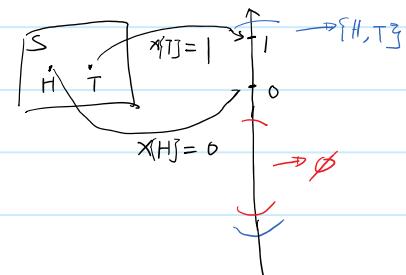
$P(\text{Event}) = \int \text{pdf}$

$P(\text{Event}) = \int \text{pdf}$

$P(\text{Event}) = \int \text{pdf}$

$$X^{-1}(A) = \{w \in S \mid X(w) \in A\}$$

inverse image $X^{-1}(A) \in \mathcal{F}$ interval



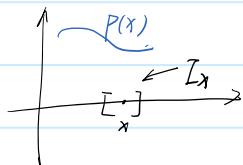
② random variable,

measurable mapping from S to \mathbb{R}

3) Pdf, (probability density function)

X: random variable, eg,

$$\overline{P}_X(x) = \overline{P}(X) = \lim_{|I_x| \rightarrow 0} \frac{P[X \in I_x]}{|I_x|}$$



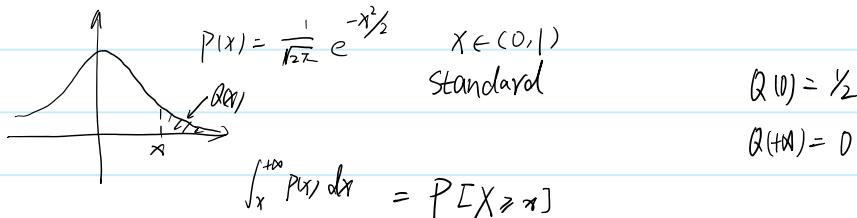
$$P[\text{event}] = P[X \in A] = \int_A p(x) dx$$

Normal

4) Gaussian r. v.

$$X_2 \quad P(X) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \text{mean } \downarrow \quad \text{Variance.} \quad X \sim N(\mu, \sigma^2) \quad N(\mu, \sigma^2)$$

\mathcal{Q} -function: tails of a Gaussian.



5) Joint & conditional Potts.

Joint p.d.f. of X and Y . e.g. (X : length., Y : width)

$$\hat{P}(x, y)$$

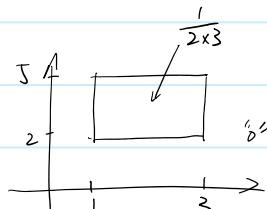
$$P \left[\text{event } C_{\text{related}} \rightarrow x \& y \right] = P \left[(x, y) \in A \right]$$

$$= \iint_A P(x, y) dx dy.$$

P [length > 2, width < 3]

$$= \int_2^3 \int_2^3 \frac{1}{6} dx dy$$

$$= \int_2^3 dx \int_2^3 dy \frac{1}{\delta} = \frac{1}{6}$$



Conditional

$$P(x|y) = \text{PDF of } X \text{ given } Y=y.$$

$$= \frac{P(x,y)}{P(y)} \times \text{marginal}$$

$$P(y) = \int_{-\infty}^{+\infty} P(x,y) dx$$

$$P(x|y) = \frac{P(y|x) P(x)}{P(y)}$$

cause
or
data not observable.

obs. data. likelihood prior normalization factor reverse influence.

$P(y) = \int_{-\infty}^{+\infty} dx P(y|x) P(x).$

Example: \star Multivariate Gaussian.

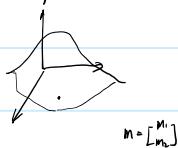
Joint Gaussian pdf.

$$x = (x_1, x_2, x_3, \dots, x_n)^T = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \vec{x}$$

$$P(x) = P(x_1, x_2, \dots, x_n) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\vec{x}-\mu)^T \Sigma^{-1} (\vec{x}-\mu)}$$

covariance matrix
 Σ
 μ mean vector
 $|\Sigma|$ determinant.

Example, $n=2$.



$$\begin{aligned} x &= \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \text{partition} \\ M &= \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \begin{array}{l} \leftarrow n_1 \text{ dim} \\ \leftarrow n_2 \text{ dim} \end{array} \\ \Sigma &= \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad \begin{array}{l} \leftarrow n_1 \times n_1 \\ \leftarrow (n-n_1) \times (n-n_2) \end{array} \end{aligned}$$

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x)}{P(x_2)}$$

$$P(x_2) = \int dx_1 P(x_1|x_2) P(x_1)$$

n-fold integration.

$$\left[\begin{array}{c} x^* \\ \Sigma \end{array} \right] \left[\begin{array}{c} x \\ \Sigma \end{array} \right] > 0 \quad \text{Positive definite}$$

Conclusion: x_1 given x_2 is Gaussian

$$P(x_1|x_2) \sim N(\mu_{1|2}, \Sigma_{1|2})$$

$$P(x_1|x_2) = 2\pi^{-\frac{1}{2}} |\Sigma_{1|2}|^{-\frac{1}{2}} \cdot e^{-\frac{1}{2}(\vec{x}-\mu_{1|2})^T \Sigma_{1|2}^{-1} (\vec{x}-\mu_{1|2})}$$

conditional mean
conditional covariance matrix

$$\mu_{1|2} = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (\vec{x}_2 - \mu_2)$$

$$\Sigma_{1|2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

non-negative revised term after condition.

measurable space

A **measurable space** is a set with a distinguished σ -algebra of subsets (called measurable). More formally, it is a pair (X, \mathcal{A}) consisting of a set X and a σ -algebra \mathcal{A} of subsets of X .

Let (X, \mathcal{A}) and (Y, \mathcal{B}) be measurable spaces.

- A map $f : X \rightarrow Y$ is called *measurable* if $f^{-1}(B) \in \mathcal{A}$ for every $B \in \mathcal{B}$.
- These two measurable spaces are called *isomorphic* if there exists a bijection $f : X \rightarrow Y$ such that f and f^{-1} are measurable (such f is called an *isomorphism*).

Let X be a set, (Y, \mathcal{B}) a measurable space, and $(f_i)_{i \in I}$ a family of maps $f_i : X \rightarrow Y$. The σ -algebra *generated* by these maps is defined as the smallest σ -algebra \mathcal{A} on X such that all f_i are measurable from (X, \mathcal{A}) to (Y, \mathcal{B}) . More generally, one may take measurable spaces (Y_i, \mathcal{B}_i) and maps $f_i : X \rightarrow Y_i$. On the other hand, if Y is \mathbb{R} (or \mathbb{C} , \mathbb{R}^n etc.) then \mathcal{B} is by default the Borel σ -algebra.

From <https://www.encyclopediaofmath.org/index.php/Measurable_space>

notation:

1. $X(\omega)$ to denote the **numerical value** of a random variable X
2. the **probability** that the outcome of the experiment is such that X is no larger than some c , i.e., that the **outcome** belongs to the **set** $\{\omega \mid X(\omega) \leq c\}$.

in order to have a probability assigned to that set, we need to make sure that it is **F-measurable**.

Definition 1. (Random variables) Let (Ω, \mathcal{F}) be a measurable space.

- (a) A function $X : \Omega \rightarrow \mathbb{R}$ is a **random variable** if the set $\{\omega \mid X(\omega) \leq c\}$ is \mathcal{F} -measurable for every $c \in \mathbb{R}$.
- (b) A function $X : \Omega \rightarrow \overline{\mathbb{R}}$ is an **extended-valued random variable** if the set $\{\omega \mid X(\omega) \leq c\}$ is \mathcal{F} -measurable for every $c \in \overline{\mathbb{R}}$.

We note here a convention that will be followed throughout: we will always use upper case letters to denote random variables and lower case letters to denote numerical values (elements of $\overline{\mathbb{R}}$). Thus, a statement such as “ $X(\omega) = x = 5$ ” means that when the outcome happens to be ω , then the realized value of the random variable is a particular number x , equal to 5.

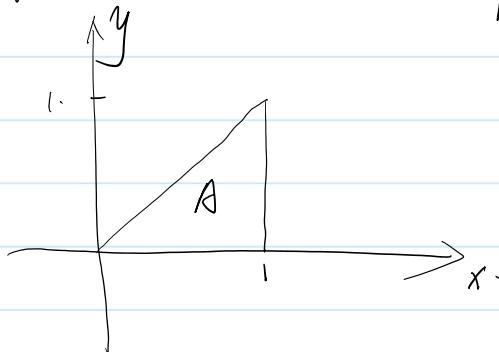
1.2 The law of a random variable

For a random variable X , the event $\{\omega \mid X(\omega) \leq c\}$ is often written as $\{X \leq c\}$, and is sometimes just called “the event that $X \leq c$.” The probability of this event is well defined, since this event belongs to \mathcal{F} . Let now B be a more general subset of the real line. We use the notation $X^{-1}(B)$ or $\{X \in B\}$ to denote the set $\{\omega \mid X(\omega) \in B\}$.

Week4

2017年9月25日 16:37

Review: Conditional PDF



$$P(x,y) = \begin{cases} c \cdot y, & (x,y) \in A \\ 0, & (x,y) \notin A \end{cases}, \quad c > 0$$

find $P(y|x) = ?$

$$P(y|x) = \frac{P(x,y)}{P(x)} = \frac{P(x \in (0,1), y \in A)}{P(x)}$$

$$\text{for } x \in (0,1) \quad P(X) = \int_{-\infty}^{+\infty} P(x,y) dy = \int_0^1 c \cdot y dy = \frac{c}{2} x^2$$

$$P(X) = \int_0^X dy c \cdot y = \frac{c}{2} x^2 \quad = \frac{c y}{x^2} \quad (x,y) \in A$$

$$= 0 \quad (x,y) \notin A$$

Random Gaussian Vector: n-dim.

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad \begin{array}{l} n_1 \text{ dim} \\ n_2 \text{ dim} \end{array}$$

$$n_1 + n_2 = n$$

$$X \sim N(m, \Sigma) \quad \begin{array}{l} \text{uncertainty} \\ \text{related to } X \end{array}$$

$$\Rightarrow m = \begin{bmatrix} m_1 \\ m_2 \end{bmatrix} \quad \begin{array}{l} n_1 \text{ dim} \\ n_2 \text{ dim} \end{array}$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \quad \begin{array}{l} n_1 \times n_1 \\ n_2 \times n_1 \\ n_2 \times n_2 \end{array}$$

$$m_{1/2} = m_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - m_2)$$

(uncertainty of x_2)⁻¹

$$\Sigma_{1/2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

Correlation between x_1 & x_2

$$P(X_1 | X_2) \sim \text{Gaussian} \sim N(m_{1/2}, \Sigma_{12})$$

$n_1 \text{ dim}$ $n_1 \times n_1$

remaining uncertainty given X_2 .

Example.

$$X = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix} \sim N(\mu, \Sigma)$$

$$\mu = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad m_1, m_2$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \Sigma_{31} & \Sigma_{32} & \Sigma_{33} \end{bmatrix}$$

2. Symmetric
↓ smaller \Rightarrow positive definite

$$m_{1/2} = m_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - m_2)$$

$$= 0 + \underbrace{\begin{bmatrix} \frac{1}{2}, \frac{1}{3} \end{bmatrix}}_{1 \times 2} \cdot \underbrace{\begin{bmatrix} 2, \frac{1}{4} \end{bmatrix}}_{2 \times 2}^{-1} \cdot \begin{bmatrix} u_2 \\ u_3 \end{bmatrix}$$

$$=$$

$$(a, b) \cdot \begin{bmatrix} u_2 \\ u_3 \end{bmatrix}$$

$$au_2 + bu_3$$

$$\underbrace{|x|}_{1 \times 1} \quad \underbrace{|x_2|}_{2 \times 1} \quad \underbrace{|x_2|}_{2 \times 1}$$

$$\Sigma_{1/2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} = 1 - \begin{bmatrix} \frac{1}{2}, \frac{1}{3} \end{bmatrix} \cdot \begin{bmatrix} 2, \frac{1}{4} \end{bmatrix}^{-1} \times \begin{bmatrix} \frac{1}{2} \\ \frac{1}{3} \end{bmatrix}$$

3) Statistics of random variables

① Expectations (mean)

$$X \text{ is r.v. } P(x) \quad E[X] = \int_{-\infty}^{+\infty} x \cdot P(x) dx$$

$$E[g(x)] = \int_{-\infty}^{+\infty} g(x) P(x) dx$$

$$E[(X - E[X])^2] = \text{Variance of } X$$

$$E[g(x,y)] = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} g(x,y) P(x,y) dx dy$$

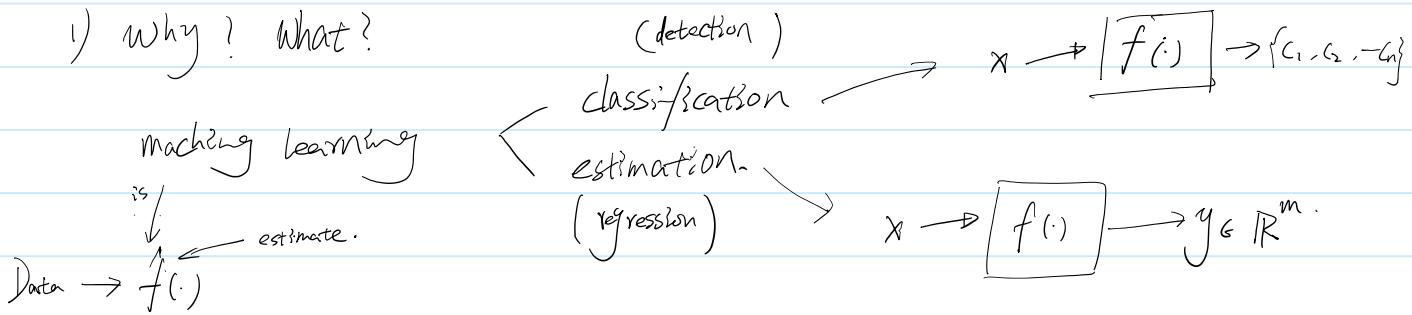
Special case.

$$E[XY] = \text{Correlation}$$

$$E[(X - E[X])(Y - E[Y])] = \text{Cov}[X, Y]$$

2. Classical theory of classification & estimation.

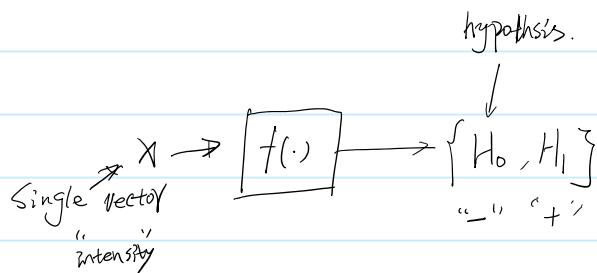
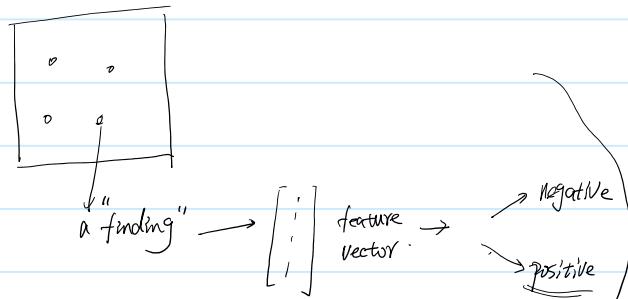
1) why? what?



classical theory, what's the best $f(\cdot)$.
(with ideal assumption)

2) classical detection theory.

① Simple example: tumor detection



Question: What's the best possible $f(\cdot)$.

② modeling,

- $P(H_0) = P_0 \quad P(H_1) = P_1$
 $P_0 + P_1 = 1.$
- $P(X|H_0)$ = conditional pdf given H_0
- $P(X|H_1)$ = conditional pdf given H_1

③ optimal $f(\cdot)$: (classifier)

Quality measure: Prob of error P_e

Consider 3 "doctors"

doctor #1, $f(x) = H_0$ optimist.
"negative"

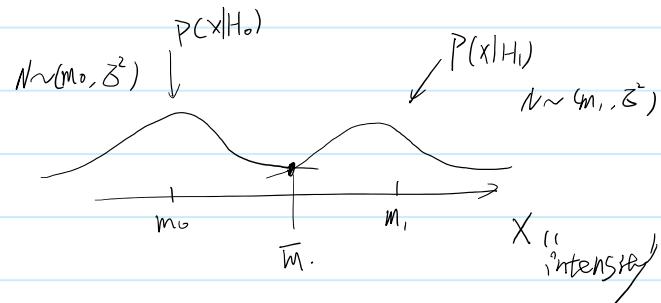
doctor #2, $f(x) = H_1$, pessimist
"positive"

doctor #3, $f(x) = \begin{cases} H_0, & \text{if } x < \bar{m}. \\ H_1, & \text{if } x \geq \bar{m}. \end{cases}$, smart

$$\bar{m} - m_0 = \frac{1}{2}(m_1 + m_0) - m_0 \\ = \frac{1}{2}(m_1 - m_0)$$

$$Q(x) = \int_x^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$$

(tail function).



Calculate P_e

$$P_e = P_{e,0} \cdot P_0 + P_{e,1} \cdot P_1 \quad \left| \begin{array}{l} \text{Given } H_0, P_{e,0} = 0 \\ \text{Given } H_1, P_{e,1} = 1 \end{array} \right.$$

P_e for each doctor,

$$\#1. \quad P_{e,0} = P[\text{error} | x \text{ comes from } H_0] = 0 \quad H_0 \text{ is true}$$

$$P_{e,1} = P[\text{error} | x \text{ comes from } H_1] = 1.$$

$$P_e = 0 \times P_0 + 1 \times P_1 = P_1$$

$$P_e = 1 \times P_0 + 0 \times P_1 = P_0$$

$$P_e = P_0 \times \int_{\bar{m}}^{+\infty} P(X|H_0) dx + P_1 \int_{-\infty}^{\bar{m}} P(X|H_1) dx$$

$$P_e = Q\left(\frac{\bar{m} - m_0}{\sigma}\right) \rightarrow \sqrt{SNR}$$

Summary:

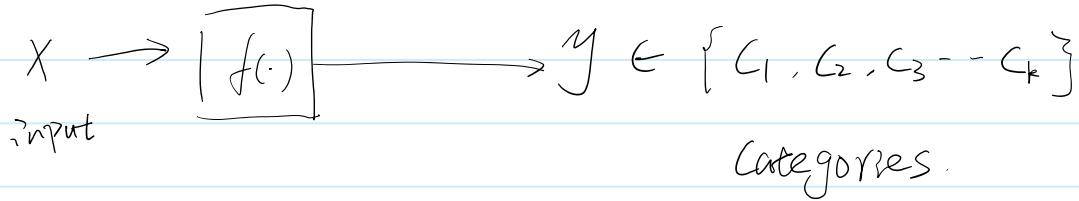
\Rightarrow None are optimal

$\min P_e$ for all $P_0, P(X|H_0), P(X|H_1)$

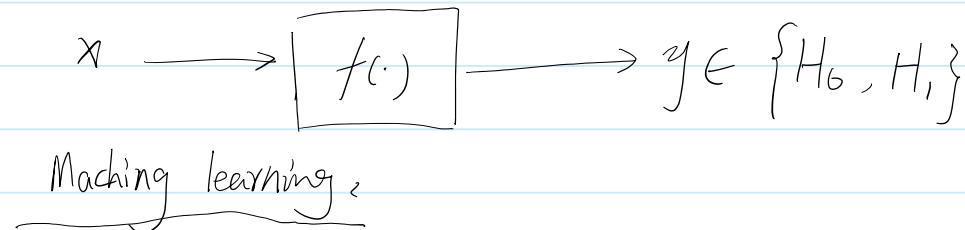
Week5

2017年10月2日 16:37

Review, classical theory of classification
Detection
Recognition



Special case (Binary)



$D, \{(x, y)\} = \text{training data.}$

$D \rightarrow f(\cdot) \text{ good classifier}$

classical theory:

model of $x \& y$ \rightarrow best possible $f(\cdot)$

$$P(H_0) = P_0 \quad (P_1 = 1 - P_0)$$

$$P(x|H_0), \quad P(x|H_1)$$

$\begin{matrix} \uparrow \\ y=H_0 \end{matrix} \quad \begin{matrix} \uparrow \\ y=H_1 \end{matrix}$

model

$\rightarrow f(\cdot)$ best "in some way" ($\min P_e$)

② modeling,

- $P(H_0) = P_0 \quad P(H_1) = P_1$
 $P_0 + P_1 = 1.$
- $P(X|H_0) = \text{conditional pdf given } H_0$
- $P(X|H_1) = \text{conditional pdf given } H_1$

③ optimal $f(\cdot)$: (classifier)

Quality measure: Prob of error P_e

Consider 3 "doctors"

doctor #1, $f(x) = H_0$ optimist.
"negative"

doctor #2, $f(x) = H_1$, pessimist
"positive"

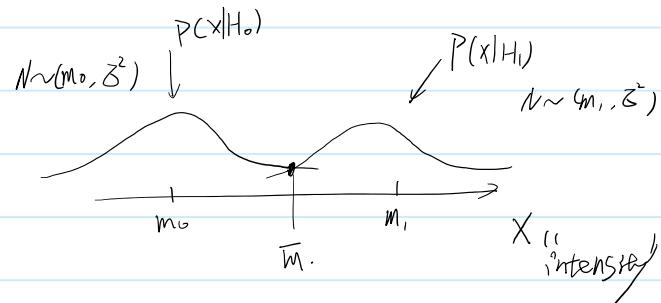
doctor #3, $f(x) = \begin{cases} H_0, & \text{if } x < \bar{m}, \\ H_1, & \text{if } x \geq \bar{m}. \end{cases}$, smart

$$\bar{m} - m_0 = \frac{1}{2}(m_1 + m_0) - m_0 \\ = \frac{1}{2}(m_1 - m_0)$$

Δm

$$Q(x) = \int_x^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{u^2}{2}} du$$

(tail function).



Calculate P_e

$$P_e = P_{e,0} \cdot P_0 + P_{e,1} \cdot P_1 \quad \left| \begin{array}{l} \text{Given } H_0, P_{e,0} = 0 \\ \text{Given } H_1, P_{e,1} = 1 \end{array} \right.$$

P_e for each doctor,

$$\#1. \quad P_{e,0} = P[\text{error} | x \text{ comes from } H_0] = 0 \quad H_0 \text{ is true}$$

$$P_{e,1} = P[\text{error} | x \text{ comes from } H_1] = 1.$$

$$P_e = 0 \times P_0 + 1 \times P_1 = P_1$$

$$P_e = 1 \times P_0 + 0 \times P_1 = P_0$$

$$P_e = P_0 \times \int_{\bar{m}}^{+\infty} P(X|H_0) dx + P_1 \int_{-\infty}^{\bar{m}} P(X|H_1) dx$$

$$P_e = Q\left(\frac{\bar{m} - m_0}{\sigma}\right) \rightarrow \sqrt{SNR}$$

Summary:

\Rightarrow None are optimal

$\min P_e$ for all $P_0, P(X|H_0), P(X|H_1)$

Optimal classifier derivation..

$$\textcircled{1} \text{ Any } f(\cdot) \equiv (R_0, R_1)$$

$$R_1 \cup R_0 = \mathbb{R}$$

$$\textcircled{2} \quad P_e = P_{e|p} P_0 + P_{e|1} P_1$$

$$\uparrow \\ P[x \in R_1 | H_0]$$

$$\Rightarrow f^*(\cdot) \text{ (optimal } f(\cdot) \text{)}$$

$$P_e = P[x \in R_1 | H_0] \cdot P_0 + P[x \in R_0 | H_1] \cdot P_1$$

$$= P_0 \int_{R_1} P(x | H_0) + P_1 \int_{R_0} P(x | H_1)$$

$$= P_0 \int_{R_1} P(x | H_0) + P_1 \left[1 - \int_{R_1} P(x | H_1) \right]$$

$$= P_1 \int_{R_1} \left(P_0 \cdot P(x | H_0) - P_1 \cdot P(x | H_1) \right)$$

$$= \underline{P_e(R_1)} - \text{a function of } R_1$$

$$\underline{P_1 P_1 P(x | H_1) > P_0 P(x | H_0)} \quad \} \quad \Leftrightarrow \text{a set of } \underline{R_1} \text{ (interval range)}$$

$\Rightarrow P_e$ minimized.

$$R_0^* = (R_1^*)^c$$

subset of real line.

Optimal classifier \Rightarrow leads to R_1

$$P_1 P(x | H_1) \underset{H_0}{\gtrsim} P_0 P(x | H_0)$$

$$\begin{array}{c} \text{likelihood} \\ \text{ratio} \end{array} \quad \frac{P(x | H_1)}{P(x | H_0)} \underset{H_0}{\gtrsim} \frac{P_1}{P_0}$$

optimal classifier
Bayesian classifier

$$\text{OR} \\ \log \left(\frac{P(x | H_1)}{P(x | H_0)} \right) \underset{H_0}{\gtrsim} \log \left(\frac{P_1}{P_0} \right)$$

$$\text{Example } P_0 \quad (P_1 = 1 - P_0). \quad P(x|H_0) \sim N(m_0, \sigma^2) \\ P(x|H_1) \sim N(m_1, \sigma^2)$$

Optimal classifier.

$$\log \frac{P(x|H_1)}{P(x|H_0)} \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} \frac{P_0}{P_1}$$

$$\begin{aligned} & \log \frac{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-m_1)^2}{2\sigma^2}}}{\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-m_0)^2}{2\sigma^2}}} \\ &= -\frac{(x-m_1)^2}{2\sigma^2} + \frac{(x-m_0)^2}{2\sigma^2} \\ &= \frac{1}{2\sigma^2} [(x-m_0)^2 - (x-m_1)^2] = \frac{1}{2\sigma^2} \left[\cancel{x} - \frac{(m_0+m_1)}{2} \right] (m_0 - m_1) \\ &= \frac{1}{\sigma^2} (x - \bar{m}) \Delta m. \end{aligned}$$

$$\begin{aligned} & \Rightarrow \frac{1}{\sigma^2} (x - \bar{m}) \Delta m \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} \log \left(\frac{P_0}{P_1} \right) \\ & (x - \bar{m}) \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} \frac{\sigma^2}{\Delta m} \log \left(\frac{P_0}{P_1} \right). \\ & x \stackrel{H_0}{\leq} \frac{\sigma^2}{\Delta m} \log \left(\frac{P_0}{P_1} \right) + \bar{m}. \end{aligned}$$

$\overbrace{\hspace{10em}}$

$P_0 = P_1$
or
 Δm is very large

Compare:

doctor #1

$$x \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} +\infty$$

doctor #2

$$x \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} -\infty$$

doctor #3

$$x \stackrel{H_1}{\geq} \stackrel{H_0}{\leq} \bar{m}.$$

Extension to more than two classes.

$$\text{Model} \quad P_k = P(H_k) \quad k = 0, 1, 2, \dots, k-1 \quad (k>2) \\ P(x|H_k)$$

Optimal classifier. Assign x to H_k^*

$$\text{if } k^* = \arg \max_{0 \leq k \leq k-1} P(H_k|x) = \arg \max_{0 \leq k \leq k-1} \frac{P(x|H_k) P_k}{P(x)} = \arg \max_{0 \leq k \leq k-1} (\log P(x|H_k) + \log P_k)$$

$$x^* = \arg \min_x (x-1)^2 \\ x^* = 1, \quad (x^* - 1)^2 = 0$$

Extension to multiple dimensions

$$P(X|H_k) \rightarrow P(X|H_k)$$

↑ ↑
1-dim n-dim

Example,

$$H_0, \quad X \sim N(m_0, \Sigma_0)$$

$$H_1, \quad X \sim N(m_1, \Sigma_1)$$

$$m_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad m_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$$\Sigma_0 = I = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \Sigma_1 = I$$

Find optimal classifier, when $P_0 = P_1 = \frac{1}{2}$.

Region, R_0 and R_1 ,

n-dim Gaussian $N(m, \Sigma)$

$$P(x) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(x-m)^T \Sigma^{-1} (x-m)}$$

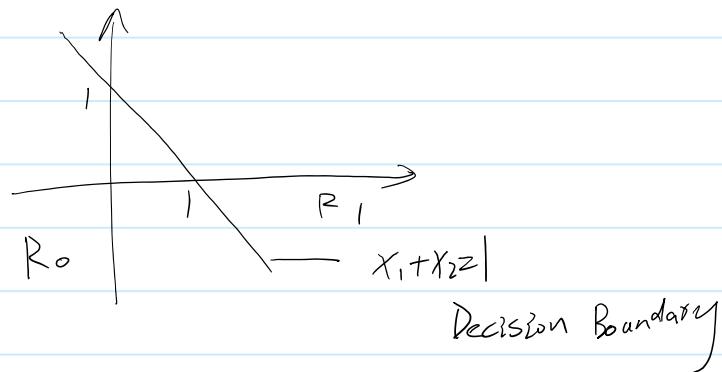


$$m \rightarrow m_1$$

$$n \rightarrow m_0$$

$$\log \frac{P(x|H_1)}{P(x|H_0)} \stackrel{H_1}{\geq} \left(\log \frac{P_0}{P_1} \right)$$

$$\text{ans} \quad X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, \quad x_1 + x_2 - 1 \stackrel{H_1}{\geq} 0 \quad \stackrel{H_0}{<} 0$$



Week6

2017年10月9日 16:37

Review,

Optimal (Bayesian classifier)

X : input

y : output, $y \in \{H_0, H_1\}$

models, $P_0 = P[H_0]$ $P_1 = P[H_1]$

$P(X|H_0)$, $P(X|H_1)$

Optimal classifier, LRT

$$\log \frac{P(X|H_1)}{P(X|H_0)} \stackrel{H_1}{\gtrless} \log \frac{P_0}{P_1}$$

↑
threshold.

General LRT

$$\log \frac{P(X|H_1)}{P(X|H_0)} \stackrel{H_1}{\gtrless} \stackrel{H_0}{\lambda} \downarrow$$

by λ .

when $\lambda = \frac{P_0}{P_1}$, Bayesian classifier

$P_{FA} = \text{Prob of false alarm} = P_{FA}$

$P_{miss} = \text{Prob of miss} = P_m$

$P_D = 1 - P_m$

Prob of detection.

defection
theory

$$\begin{cases} P_D = P_D(\lambda) \\ P_{FA} = P_{FA}(\lambda) \end{cases} \Rightarrow \text{ROC curve. } (P_{FA}(\lambda), P_D(\lambda))$$

Receiver operating characteristics.

Example:

$$P(X|H_0) \sim N(-1, 1)$$

$$P(X|H_1) \sim N(1, 1)$$

$$\log \frac{P(X|H_1)}{P(X|H_0)} \geq \lambda$$

$$\text{assume } P_0 = P_1 = \frac{1}{2} \quad P_D = P_D(\lambda) = \dots$$

find P_D vs P_{FA} as a function of λ .

$$P_{FA} = P_{FA}(\lambda) = \dots$$

Practical problems with Optimal classifier

Training data. $D = \{(x_i, y_i)\}_{i=1}^n$

\downarrow \uparrow
 real from \mathbb{R}^m . H_0 or H_1

Don't have, $P(X|H_0)$, $P(X|H_1)$.

Solution.

① density estimation.

$D \rightarrow \hat{P}(x|H_0) \cdot \hat{P}(x|H_1)$: This requires simple model
no work well, with high-D data : ex. non-Gaussian

② $D \rightarrow$ classification function $f(\cdot)$.

3) Classical Estimation Theory

① Estimation problems.

Density estimation

Variable / function estimation

Density estimation.

$$D = \{x_1, x_2, \dots, x_n\}$$

= i.i.d. Samples of X

independent, identically distributed.

$$D \longrightarrow \hat{P}(x)$$

Density estimation.

e.g. Gaussian

✓ Parametric : $P(x) = P(x|\theta)$ $\theta \sim (\mu, \sigma^2)$

Non-Parametric

e.g. MLE

✓ frequentist : no assumptions on $p(x)$ "complete data driven"

✓ Bayesian : assumes some prior on $p(x)$

e.g. $P(\theta)$ MAP

Variable / function estimation.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

= random, iid. Samples of X and Y .
(X, Y).

$$D \rightarrow \hat{Y} = f(X)$$

② Parametric density estimation.

(parameters estimation)

frequentist : MLE (max likelihood estimation)

Bayesian : MAP (max a posteriori)

Example: Estimating a mean of a Gaussian.

$$X_1 \ X_2$$

$D = \{x_1, x_2, \dots, x_n\} = \text{i.i.d. samples of } X$

$$X \sim N(m, \sigma^2)$$

σ^2 known.

m unknown - to be estimated from D

$$P(X|\theta) = P(X|m) \quad \theta = m$$

MLE,

$$\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} P(x_1, \dots, x_n | \theta) \quad \text{: always exponential}$$

$$= \arg \max_{\theta} \log P(x_1, \dots, x_n | \theta) \quad \text{take log}$$

How good is \hat{m} .

$$\rightarrow E[\hat{m}] = ?$$

$$\underset{\text{Unbiased.}}{=} m.$$

$$(a+b+c-\bar{z})^2 = a^2 + b^2 + c^2 - 2ab - 2ac - 2bc + 2\bar{z}a + 2\bar{z}b + 2\bar{z}c - 3\bar{z}^2$$

$$= 2ab + \underbrace{2ac + 2bc}_{C_{26}} - 3\bar{z}^2$$

$$\rightarrow E[(\hat{m} - m)^2] = \text{MSE} = \text{Var}[\hat{m}]$$

↑ true unknown parameter.

$$\begin{aligned} E[\hat{m}] &= E\left[\frac{1}{n} \sum x_i\right] \\ &= \frac{1}{n} \sum E[x_i] \\ &= \frac{1}{n} \sum m = \frac{nm}{n} = m \end{aligned}$$

$$E[(\hat{m} - m)^2] = E\left[\left(\frac{1}{n} \sum x_i - m\right)^2\right]$$

$$\begin{aligned} &= E\left[\left(\frac{1}{n} \sum x_i\right)^2 - 2m \frac{1}{n} \sum x_i + m^2\right] \\ &= \frac{1}{n} \sum x_i^2 - 2m \frac{1}{n} \sum x_i + m^2 \\ &= -2\bar{m}m + m^2. \end{aligned}$$

=

5

2017年10月9日 16:37

Q function and Error functions

2017年10月9日 16:37

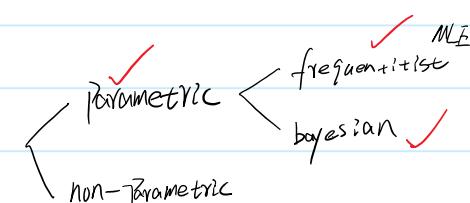
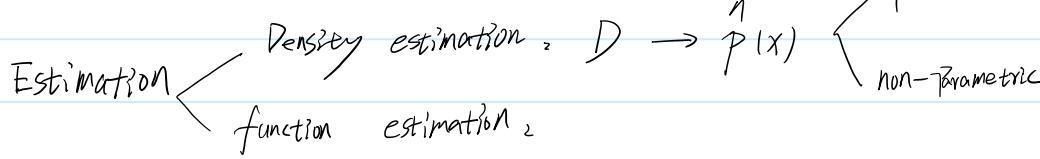
<https://www.gaussianwaves.com/2012/07/q-function-and-error-functions/>

Week7

2017年10月16日

16:37

Review, MLE



MLE under certain condition

become the real θ when $|D| \rightarrow \infty$

$$\text{MLE: } p(x) = P(x|\theta)$$

$$\hat{\theta} = \arg \max_{\theta} P(D|\theta)$$

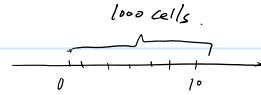
$$= \arg \max_{\theta} \log P(D|\theta)$$

$$E[(\hat{\theta}(n) - \theta)^2] \xrightarrow{n \rightarrow \infty} 0$$

↓
sample size

X : 1 dim.

100,000 points.



enough data

X : 1000 dim.

100,000 points

cells: 1000^{100}

not enough data.

Bayesian estimation.

can be useful when, don't have enough data.

MAP (max a posterior)

Conditional mean

$P(X) = P(X|\theta)$, likelihood function.

$P(\theta)$ = prior density for θ .

"knowledge/constraint" about θ

MAP:

$$\hat{\theta}_{\text{map}} = \arg \max_{\theta} P(\theta|D) = \arg \max_{\theta} P(D|\theta) P(\theta) / P(D)$$

$$\begin{aligned} \text{MAP: } \hat{\theta}_{\text{map}} &= \arg \max_{\theta} \{ P(D|\theta) P(\theta) \} \\ &= \arg \max_{\theta} \{ \log P(D|\theta) + \log P(\theta) \} \end{aligned}$$

if $\log P(\theta) = \text{Constant}$.
then MAP \Rightarrow MLE
means "no knowledge about θ "

Example 2 Gaussian Case.

$$\begin{aligned} D &= \{x_1, x_2, \dots, x_n\} \\ &= \text{i.i.d. samples of } X \end{aligned}$$

$$X \sim N(m, \sigma^2)$$

Unknown Known

$$\theta = m$$

$$P(\theta) = P(m) \sim N(m_0, \sigma_0^2)$$

assumed to be known.

$$\hat{\theta}_{\text{map}} = ?$$

③ Variable (function estimation)

Recall, the problem.

$$\begin{aligned} D &= \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\} \\ &= \text{i.i.d. samples of } X \text{ and } Y \end{aligned}$$

estimation: $D \rightarrow \hat{f}(.)$ $\hat{Y} = \hat{f}(x)$

$\hat{f}(.)$: estimation of $f(.) = Y = f(X)$

Theoretical result,

$$\hat{Y} = g(X) \quad \text{any estimate}$$

optimal estimate in sense of MMSE.

$$g^*(X) = \mathbb{E}[Y|X]$$

$\hat{Y} = \mathbb{E}[Y|X]$ is the optimal estimator of Y from X .

Week8

2017年10月23日

16:37

Review: Variable / function estimation

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

= i.i.d. subsets from (X, Y)

estimation, $D \rightarrow f(\cdot) \Rightarrow \hat{Y} = f(x) \approx Y$

base estimate

$$\hat{Y} = \hat{f}(x) = E[Y|X]$$

= best estimate for Y
in MSE sense.

"Proof." Mean Squared Error

$$\begin{aligned} \mathcal{E} &= E[(y - \hat{y})^2] \\ &= \iint (y - \hat{y})^2 P(x, y) dx dy \\ &= \int dx \int dy (y - \hat{y})^2 P(x) P(y|x) \\ &= \int dx P(x) \int dy (y - \hat{y})^2 P(y|x) \end{aligned}$$

$\overbrace{\quad\quad\quad}^{\substack{E[h(x,y)] \\ = \int h(x,y) p_{x,y} dx dy}}$
 $\overbrace{\quad\quad\quad}^{\substack{\text{joint prob} \\ P(x,y) \\ = P(x)/P(y|x) \\ = P(y) P(x|y)}}$

$$\Rightarrow \text{To min } \mathcal{E} \Rightarrow \min \int (y - \hat{y})^2 \text{ mean squared error}$$

$P(y|x) dy$ for each x .

$\sum a_i, a_i \geq 0$
 $\min \text{ each } a_i$.

for each x , min $\int dy (y - \hat{y})^2 P(y|x)$

$$\frac{\partial}{\partial \hat{y}} \int dy (y - \hat{y})^2 P(y|x) = 0$$

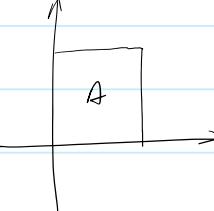
$$-\int dy 2(y - \hat{y}) P(y|x) = 0$$

$$-\int dy (y - \hat{y}) P(y|x) = 0$$

$$\begin{aligned} \int dy (y - \hat{y}) P(y|x) &= \int y P(y|x) dy - \int \hat{y} P(y|x) dy \\ \hat{y} \int dy P(y|x) &= \int y P(y|x) dy \\ \hat{y} &= \int y P(y|x) dy = E[y|x] \end{aligned}$$

$$E[(y - E[y|x])g(x)] = 0 \quad \text{orthogonality principle.}$$

Example 1. $P_{xy}(x,y) = \begin{cases} C(x+y), & (x,y) \in A \\ 0, & (x,y) \notin A \end{cases}$



find the optimal estimate of y from x .
 $= E[y|x]$

Example 2.

In practice:

$$D \rightarrow \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$D \rightarrow f(x) \approx E[y|x]$$

① $D \rightarrow \hat{P}(x,y) \rightarrow E[y|n] \quad | \text{ difficult when } \dim(x), \dim(y) \text{ are high.}$
density estimation

\rightarrow assume $P(x,y)$ is Gaussian $|$ Problem: Data is not Gaussian

② find/start with a family of functions $y = g(x) = g(x|\theta) \leftarrow \text{parameters.}$
 $= \hat{g}(x|\theta)$

$$\hat{\theta} = \arg \min \sum (y_i - g(x_i|\theta))^2 / n$$

$E[(y - \hat{g})^2]$

If $g(x|\theta)$ is rich

$$\Rightarrow g(x|\hat{\theta}) \approx E[y|x]$$

$$\frac{\sum y_i}{n} = \text{Sample mean} \rightarrow E[\bar{Y}]$$

$n \rightarrow \infty$

4. Linear and logistic regression

1) Linear regression

function estimation,

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

find: $f(\cdot)$ = linear.

$$\hat{y} = w^T x + b$$

↑ ↑ ↑
 Scalar Row Vector Scalar
 ↓ ↓ ↓
 Column Vector

Example, $x \in 2\text{dim}$

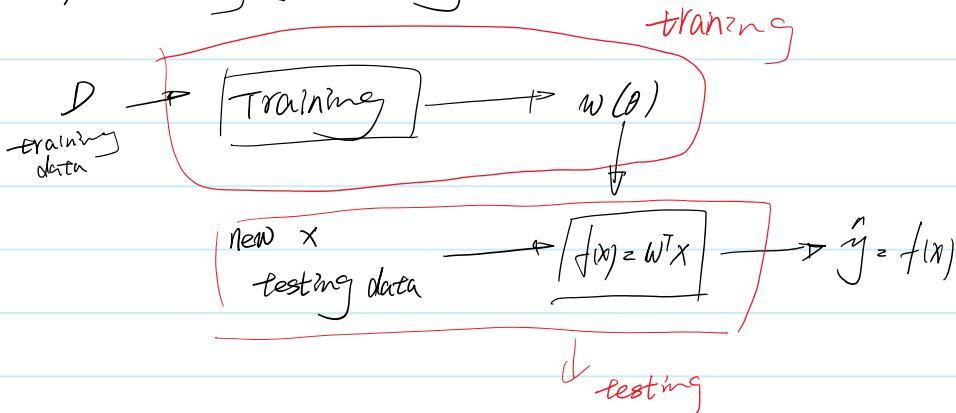
$$y: 1\text{dim} \quad f(x) = w_1 x_1 + w_2 x_2 + w_0$$

$$= [w_1, w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + w_0$$

$$= [w, w_2, w_0] \begin{bmatrix} x_1 \\ x_2 \\ 1 \end{bmatrix}$$

$$\theta = [w, w_2, w_0]^T$$

2) Training (and testing)



"Training" $D \rightarrow w$

"Training" $D \rightarrow w$

① 1st perspective / derivation, data-driven

least squared

$$\text{recall } \mathcal{E} = \text{MSE} = \frac{\sum (y_i - f(x_i))^2}{n}. \quad f(x) = w^T x \quad \xleftarrow{\text{augmented}} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\underset{\text{optimal}}{\downarrow} \quad = \mathcal{E}(w)$$

$$w^* = \arg \min_w \mathcal{E}(w)$$

$$\nabla_w \mathcal{E} = 0 \Rightarrow w^*$$

$$F(x) = F(x_1, x_2, \dots, x_n)$$

$$\frac{\partial \mathcal{E}}{\partial w} = 0$$

$$\nabla_{x_i} F = \frac{\partial F}{\partial x_i} \quad i = \text{gradient}$$

$$\boxed{\nabla_w \mathcal{E} = 0}$$

normal equation

$$\frac{\partial F}{\partial x_n}$$

Example: $x_1 = 1d$

$y_1 = 1d$.

$$y = f(x) = w_1 x_1 + w_0 = [w_1, w_0] \begin{bmatrix} x_1 \\ 1 \end{bmatrix}$$

$$\mathcal{E} = \frac{1}{n} \sum (y_i - (w_1 x_i + w_0))^2$$

$$\nabla_w \mathcal{E} = 0, \quad \frac{\partial \mathcal{E}}{\partial w_1} = 0, \quad \frac{\partial \mathcal{E}}{\partial w_0} = 0$$

$$\frac{\partial \mathcal{E}}{\partial w_1} = \frac{1}{n} \sum 2 \cdot (y_i - (w_1 x_i + w_0)) \cdot (-x_i) = 0 \Rightarrow (\sum x_i^2) w_1 + (\sum x_i) w_0 = \sum x_i y_i$$

$$\frac{\partial \mathcal{E}}{\partial w_0} = \frac{1}{n} \sum (y_i - (w_1 x_i + w_0)) \cdot 1 = 0 \Rightarrow (\sum x_i) \cdot w_1 + (\sum 1) \cdot w_0 = \sum y_i$$

$$\begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & n \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

$$A \cdot w = b$$

$$w^* = A^{-1} \cdot b$$

Another way.

$$f(x) = w_1 x + w_0$$

$$\mathbb{E}[\varepsilon] = \mathbb{E}[(y - f(x))^2]$$

$$\frac{\partial \varepsilon}{\partial w_1} = 0$$

$$- \mathbb{E}[2(y - (w_1 x + w_0)) \cdot x] = 0$$

$$\mathbb{E}[x^2] w_1 + \mathbb{E}[x] \cdot w_0 = \mathbb{E}[xy]$$

$$\frac{\partial \varepsilon}{\partial w_0} = 0$$

$$- \mathbb{E}[2(y - (w_1 x + w_0)) \cdot 1] = 0$$

$$\mathbb{E}[x] w_1 + 1 \cdot w_0 = \mathbb{E}[y]$$

$$\begin{bmatrix} \mathbb{E}[x^2] & \mathbb{E}[x] \\ \mathbb{E}[x] & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} = \begin{bmatrix} \mathbb{E}[xy] \\ \mathbb{E}[y] \end{bmatrix}$$

A

$$\begin{bmatrix} \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & n \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n x_i y_i \\ \sum_{i=1}^n y_i \end{bmatrix}$$

General Case x_i : m -dim

y_i : 1-dim.

$$\mathcal{E} = \frac{1}{n} \sum_i (y_i - x_i^T w)^2 \quad \text{← } m+1 \text{ dim.}$$

$$\nabla_w \mathcal{E} = 0$$

$$\hookrightarrow (X^T X)w = X^T y$$

$$X = \begin{pmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{pmatrix} \quad n \times M$$

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

$$w = (X^T X)^{-1} X^T y$$

$\begin{matrix} \checkmark & \downarrow & \downarrow & \downarrow \\ m \times n & n \times m & m \times n & n \times 1 \\ \downarrow & & \swarrow & \\ m \times m & & m \times n - n \times 1 \Rightarrow m \times 1 & \end{matrix}$

Week9

2017年10月30日

20:18

Review : Linear regression

$$D = \{(x, y)\}$$

$D \rightarrow w$ (training)

$$D \rightarrow f()$$

$$f(\cdot) : x \rightarrow y$$

$$\hat{y} = f(x) = w^T x \quad \text{augmented}$$

① least squares (data driven)

② random variable based.

③ Gaussian based.

Example: x_i : 1d samples $\rightarrow x$

y_i : 1d samples $\rightarrow y$

$$w = \begin{bmatrix} w_1 \\ w_0 \end{bmatrix}$$

$$\text{find } w \text{ : } \min \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$\begin{aligned} \hat{y} &= w_1 x_1 + w_0 \\ &= \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} [x_1, 1] \end{aligned}$$

$$\begin{bmatrix} \sum x_i^2 & \sum x_i \\ \sum x_i & n \end{bmatrix} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} = \begin{bmatrix} \sum x_i y_i \\ \sum y_i \end{bmatrix}$$

Assumes:

$$x_i \sim X$$

$$y_i \sim Y \quad \text{② min } E[(y - w^T x)^2]$$

$$E[X] = \frac{\sum x_i}{n}$$

$$\Rightarrow \begin{bmatrix} E[x^2] & E[x] \\ E[x] & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_0 \end{bmatrix} = \begin{bmatrix} E[xy] \\ E[y] \end{bmatrix}$$

Normal
Eqn

Assumes:

$$x, y \text{ are Gaussians} \quad \text{③ } \begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} m_x \\ m_y \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}\right) \Rightarrow m_{yx} = E[y/x] \\ = m_y + \sigma_{yx} (\sigma_x^2)^{-1} (x - m_x)$$

$$m_x = E[x], m_y = E[y]$$

$$\sigma_x^2 = E[(x - E[x])^2]$$

Variance of estimation error.

$$\hat{y} = f(x)$$

$$\sigma_y^2 = E[(y - E[y])^2]$$

$$\sigma_{yx}^2 = \sigma_y^2 - \sigma_{yx} (\sigma_x^2)^{-1} \sigma_{xy}$$

$$= \underbrace{\sigma_{yx} (\sigma_x^2)^{-1}}_{w_1} x + \underbrace{m_y - \sigma_{yx} (\sigma_x^2)^{-1} m_x}_{w_0}$$

$$\sigma_{xy} = \sigma_{yx} = E[(y - E[y])(x - E[x])]$$

$$\downarrow E[(y - E[y/x])^2]$$

General Case x_i : m-dim Samples from x .
 y_i : 1-dim. Samples from y

$$\text{MSE} = \mathcal{E} = E[(y - w^T x)^2] \quad \xrightarrow{\text{augmented.}} \quad X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_m \\ w_0 \end{bmatrix}$$

$$\mathcal{E} = \frac{1}{n} \sum (y_i - x_i^T w)^2$$

$$D_w \mathcal{E} = 0$$

$$D_w \mathcal{E} = E[-2(y - w^T x) D_w(w^T x)]$$

$$= -2 E[(y - w^T x) \cdot x] = 0$$

Gradient $\xrightarrow{\text{variable vector}}$

$$D_x(a^T x)$$

\nwarrow constant vector

$$= D_x(\mathcal{E} a^T x) = \begin{bmatrix} \frac{\partial}{\partial x_1} \sum a_i x_i \\ \frac{\partial}{\partial x_2} \sum a_i x_i \\ \vdots \\ \frac{\partial}{\partial x_n} \sum a_i x_i \end{bmatrix}$$

$$= \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = a.$$

$$D_x(a^T x) = a$$

$$\underline{D_x(x^T a) = a.}$$

$$\Rightarrow E[(x^T w) \cdot x] = E[x^T y]$$

$$E[x \cdot (x^T w)] = E[x^T y]$$

$$(m+1) \times 1 \quad 1 \times (m+1)$$

$$w E[xx^T] = E[x^T y]$$

$$w = (E[xx^T])^{-1} E[x^T y]$$

Approximation (using D)

$$E[x^T y] = \frac{1}{n} \sum x_i y_i$$

$$E[xx^T] = \frac{1}{n} \sum x_i x_i^T \quad \xrightarrow{\text{augmented.}}$$

$$\begin{bmatrix} I & x \\ x^T & \end{bmatrix} \Rightarrow \begin{bmatrix} I & x \\ x^T & \end{bmatrix}$$

$$X = \begin{bmatrix} X_1^T \\ X_2^T \\ \vdots \\ X_n^T \end{bmatrix} = \begin{bmatrix} [x_1, -1] \\ [x_2, -1] \\ \vdots \\ [x_n, -1] \end{bmatrix} \stackrel{m+1}{\text{mtl.}}$$

$$X^T = \begin{bmatrix} [1] \\ [x_1] \\ [x_2] \\ \vdots \\ [x_n] \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\frac{1}{n} \sum X_i X_i^T = X^T X$$

$$X, X^T \Rightarrow = [x_1, x_2, \dots, x_n] \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_n^T \end{bmatrix} = \sum X_i X_i^T.$$

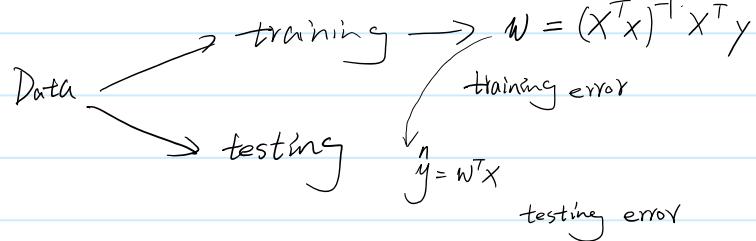
$$\sum X_i Y_i = [x_1, x_2, \dots, x_n] \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = X^T Y$$

$$\frac{1}{n} (X^T X) w = \frac{1}{n} X^T Y$$

$$w = (X^T X)^{-1} (X^T Y)$$

3) Regularized linear regression

① why?



Possible reasons:

$(X^T X)$ ill-conditioned

② Possible solution, "Regularization"

③ A common regularization, ℓ_2 minimization.

$$w = \arg \min_w E[(y - w^T x)^2]$$

$$\Rightarrow w = \arg \min_w [E[(y - w^T x)^2] + \lambda w^T w] \quad \|\omega\|_2^2 = \sum_{j=0}^{m+1} w_j^2$$

$f(w)$

$$f(w) = \mathcal{E}(w) + \lambda w^T w$$

$$\lambda > 0$$

$$\left| \begin{array}{c} \downarrow \\ D_\lambda \\ \downarrow \\ X^T A X \end{array} \right|$$

\$D_\lambda\$ \$m \times m\$ square & symmetrical
 \$X^T A X\$ \$m \times 1\$
 \$= \geq A X\$

MIN
Gradient

$$D_w F(w) = D_w \mathcal{E}(w) + \lambda D_w (w^T w)$$

$$= E[(y - w^T x) \cdot \lambda] + \lambda \cdot 2 I \cdot w$$

$$\Rightarrow (E[x x^T] + 2\lambda I)w = E[x y]$$

\downarrow
makes more invertable.

$$2.2) \text{ Variation, } \text{ l1 minimization } w^T w \rightarrow \|w\|_1 = \sum_{j=0}^m |w_j|$$

2.3) PCA analysis

$$x_{max} \rightarrow x_1 \rightarrow LR.$$

$$\left[\begin{array}{c} \vdots \\ \vdots \end{array} \right] \rightarrow \left[\begin{array}{c} \vdots \\ \vdots \end{array} \right]$$

Week10

2017年11月6日 20:18

4) Iterative Solution \approx Gradient descent

① why?

$$\text{LR equation: } W^* = (E[XX^T])^{-1} E[Xy]$$

$\dim(X) = m$

if m is very large

$\Rightarrow (E[XX^T])^{-1}$ computationally expensive

Solution

① PCA dim reduction.

$$X \xrightarrow{\text{PCA}} X' \quad \begin{matrix} \text{only work if most of} \\ \text{EV of } E[XX^T] \text{ is} \\ \text{small} \end{matrix}$$

e.g. $\sim 10^4$

$\sim 10^1$

② Iterative Solution

$$W_{n+1} = W_n + \frac{\Delta W_n}{\| \cdot \|} \quad \begin{matrix} \text{Gradient Descent} \\ \text{increment} \leftarrow \text{Newton Method} \\ \text{Pseudo} \\ + \text{limited memory} \end{matrix}$$

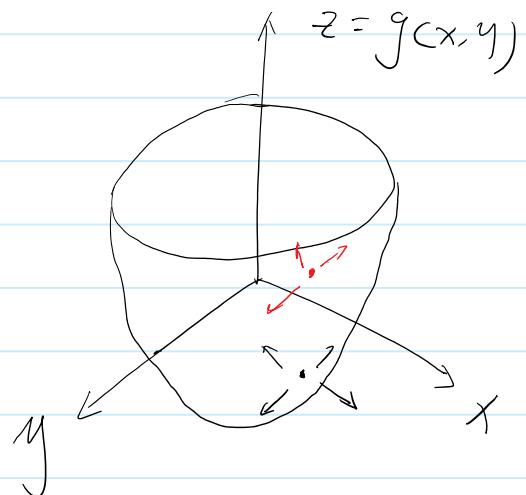
result
from previous
iteration

② Gradient Descent

Gradient
minimized.

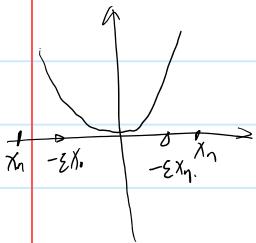
$$\Delta w_n = -\varepsilon \nabla_w E(w) \Big|_{w=w_n}$$

Recall, $E(w) = \mathbb{E}[(y - w^T x)^2]$



Example (1D)

$$g(x) = x^2$$



$$\nabla g = \frac{dg}{dx} = 2x$$

$$\begin{aligned} x_{n+1} &= x_n - \varepsilon 2x_n \\ &= x_n - \varepsilon x_n. \end{aligned}$$

$$x_{n+1} = x_n - \varepsilon x_n = (1 - \varepsilon)x_n.$$

$$= (1 - \varepsilon)(1 - \varepsilon)x_{n-1}$$

$$= (1 - \varepsilon)^n x_0$$

$$\text{if } 0 < \varepsilon < 1 \rightarrow x_n \rightarrow 0$$

$$n \rightarrow \infty$$

$$\frac{dg}{d\vec{n}} = \vec{n} \cdot \nabla g$$

$\max \left| \frac{dg}{d\vec{n}} \right|$ achieved when \vec{n} in the same/opp. direction

descent \rightarrow opp direction of ∇g .

GD $\left(\boxed{x_{n+1} = x_n - \varepsilon \nabla g(x_n)} \right)$

$$GD_2 \left(\overbrace{x_{n+1} = x_n - \varepsilon Dg(x_n)} \right)$$

applied to linear regression

$$E(w) = E[(y - w^T x)^2]$$

$$\nabla_w E(w) = -2 E[(y - w^T x) \cdot x]$$

$$= 2 \left(E[xy] - E[xx^T]w \right)$$

$$GD_2 \quad w_{k+1} = w_k - \varepsilon \left(E[xy] - E[xx^T]w_k \right)$$

in each iteration

Take a ^{random} subset of $\{x_i, y_i\}$ from entire training data.

$$E[xy] = \frac{1}{n} \sum_{i=1}^n x_i y_i \quad \text{over subset}$$

\Rightarrow min batch.

Practical implementation:

In each iteration

$$E[xy] = \frac{1}{n} \sum_{i=1}^n x_i y_i, \quad E[xx^T] = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$$

\Rightarrow Batch method

n : sample size

In each iteration:

$$\text{Take one pair } x_i, y_i, \quad E[xy] = x_i y_i, \quad E[xx^T] = x_i x_i^T$$

\Rightarrow SGD / LMS / adaptive filtering

Stochastic
Gradient
Descent

5) logistic regression

① Problem : to learn a linear classifier

$$\text{Data} \in D = \{(x_i, y_i)\}$$

x_i , n-dim

y_i 1 dim $y_i \in \{0, 1\}$

H_0, H_1

$$D \rightarrow g(x) = P[y=1]$$

$$D \rightarrow g(\cdot) \quad g(\cdot) = [g_{H_1}=1] \Rightarrow D \rightarrow w$$

\Rightarrow learning,

$$g(x) = \frac{1}{1 + e^{-w^T x}}$$

$\mathcal{E}(w)$, function to be optimised

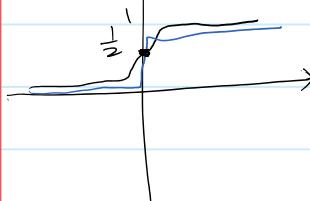
$$= \varphi(w^T x)$$

Cross entropy (non-linear non-quadratic)

$$G(D, w_{n+1} = w_n - \varepsilon \nabla_w \mathcal{E}(w)) \xrightarrow{\text{loss function}}$$

② logistic function: continuous version of threshold.

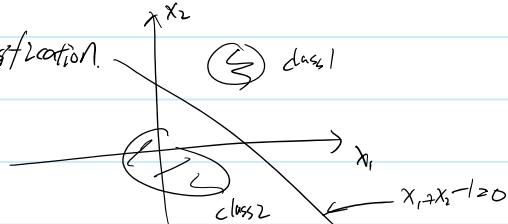
$$\varphi(t) = \frac{1}{1 + e^{-t}}$$



$$\mu(t) = \begin{cases} 0, & t < 0 \\ 1, & t \geq 0 \end{cases}$$

↓
doesn't have derivative.

relation to classification.



$\varphi(w^T x)$ generates $P[x \in H_1]$

Soft classification

Surpassed by

{ linear, SVM (kernel)
non-linear, Neural Network }

(4) Cross entropy

交叉熵

$$X \in \mathcal{P}(X)$$

$-\log_2 p(x) = \min$ # of bits needed to encode x .

X	1	2	3	4
$p(x)$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$

 $1 \rightarrow 2$ bits 00 $2 \rightarrow 2$ bits 01 Ave # of bit need = 2 $3 \rightarrow 2$ bits 10 $4 \rightarrow 2$ bits 11

X	1	2	3	4
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$

bits	1	2	3	3
0	0	10	110	111

$$\text{Ave: } 1 \times \frac{1}{2} + 2 \times \frac{1}{4} + 3 \times \frac{1}{8} + 3 \times \frac{1}{8} = \frac{1}{8}(4+4+3+3) = \frac{7}{4} < 2$$

$$\sum_x p(x) (-\log p(x)) = H(X) = \min \text{ Ave # bits need to code } X.$$

= the amount of uncertainty

= entropy

Cross entropy ↘ real
assumed.

6

2017年11月6日 20:18

7

2017年11月6日 20:18

Week11

2017年11月13日

20:18

Cross entropy:

Two
possibilities
 X_1 binary

$$P = P[X=1]$$

$$1-P = P[X=0]$$

$$g = Q[X=1]$$

$$1-g = Q[X=0]$$

$$H(P, g) = -P \log g - (1-P) \log(1-g)$$

$$H(P) = -P \log P - (1-P) \log(1-P)$$

In general

$$H(P) = -\sum P_k \log P_k$$

X	x_1, x_2, \dots, x_k
P	P_1, P_2, \dots, P_k

= min. ave. codeword length need.

$$H(P, g) = -\sum_{k=1}^k P_k \log g_k$$

= ave. codeword length by
assuming the data distribution to be g

while it actually is P

(optimize codeword length according to g)

Example:

$$P = [\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}] \quad H(P) = \frac{1}{2} + 2 \cdot \frac{1}{4} + 3 \cdot \frac{1}{8} + 3 \cdot \frac{1}{8} = \frac{7}{4} = 1.75$$

$$g = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}] \quad H(g) = 2$$

$$H(P, g) = -\sum P_k \log g_k = -[\frac{1}{2}(-2) + \frac{1}{4}(-2) + \frac{1}{8}(-2) + \frac{1}{8}(-2)]$$

$$= 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{4} = 2 > H(P) = 1.75$$

$$H(P, g) - H(P) = 2 - 1.75 = 0.25$$

$$\begin{cases} g = [\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}] \\ P = [\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}] \end{cases}$$

$$H(P, g) = \frac{9}{4} > H(P) = 2$$

$$H(P, g) - H(P) \geq 2$$

$$D(P//g) = kL \underline{\text{distance}}$$

$$D(P//g) = \sum P_k \log \frac{P_k}{g_k} \neq D(g//P)$$

Application in estimation

P : unknown distribution

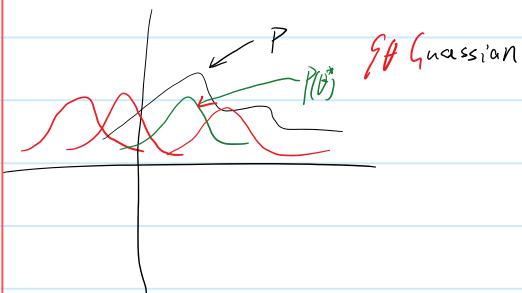
q : parametric distribution

$$q = q(\theta)$$

find an estimate of P using $q(\theta)$.

$$\theta^* = \arg \min_{\theta} D(P // q(\theta))$$

$q(\theta^*)$ best estimate of P .



$$E[X] = \int_{-\infty}^{\infty} x P(x) dx$$

$$\approx \frac{1}{n} \sum X_i$$

$$\theta^* = \arg \min_{\theta} D(P // q(\theta))$$

$$= \arg \min_{\theta} H(P, q(\theta))$$

$$D(P // q) = H(P, q) - H(P)$$

find a q that most close to P

④ Apply cross entropy \rightarrow logistic regression

recall linear regression,

$$\hat{y} = f(x) = w^T x$$

$$MSE = E[(y - w^T x)^2] \rightarrow w = E[xx^T]^{-1} E[xy]$$

logistic regression:

two class classification

x — input, m -dim

$$P(x) = P[y=1|x]$$

y — output: 1-dim $y \in \{0, 1\}$

= true distribution = unknown.

$g(x) = \text{approximation to } P(x)$

$$H(P, q)(x)$$

$$= [-P \log q - (1-P) \log(1-q)](x)$$

$$= g(w^T x) = \frac{1}{1 + e^{-w^T x}}$$

augmented

logistic regression: find $q(x)$ = best approx. to $P(x)$.

$$\Rightarrow w^* = \arg \min_w H(P, q)(x) = \arg \min_w E[H(P, q)(x)]$$

$H(p, q)(x)$

$$= [-p \log q - (1-p) \log(1-q)](x)$$

logistic regression; find $g(x) = \text{best approx. to } p(x)$

$$\Rightarrow w^* = \arg \min_w H(p, q)(x) = \arg \min_w E[H(p, q)(x)]$$

$$\Rightarrow \text{GD} \quad w_{n+1} = w_n - \varepsilon D_w \mathbb{E}[H(p, q)(x)]$$

$$\begin{cases} D_w \in [H(p, q)(x)] \\ = E[(q-p)x] \\ \approx \frac{1}{n} \sum (q(x_i) - p(x_i)) x_i \end{cases}$$

$$\approx \frac{1}{n} \sum (q(x_i) - y_i) x_i$$

$$\left\{ \begin{array}{l} P(x) = P[y=1|x] \\ y_i=1 \rightarrow P(x_i) \approx 1 \\ y_i=0 \rightarrow P(x_i)=0 \end{array} \right.$$

(5) multi class case,

input $x = m$ -dim.

output y , 1 -dim $y \in \{1, 2, \dots, k\}$

k -dim indicator vector $\vec{y} = [y_1, y_2, \dots, y_k]$

$y = k^{\text{th}}$ class = $[0, 0, \dots, 0, 1, 0, \dots, 0]^T$

$$P(x) = [P_1(x), P_2(x), \dots, P_k(x)]^T$$

$$P_k(x) = P[y = k^{\text{th}} \text{ class} | x]$$

$$q(x) = [q_1(x), q_2(x), \dots, q_k(x)]^T \quad \text{based on } g(y)$$

$$\hat{q}_k(x) = \frac{e^{w_k^T x}}{\sum_{k'} e^{w_{k'}^T x}} \quad \begin{matrix} \leftarrow & q(k) \\ \leftarrow & \text{normalization} \end{matrix} \quad \hat{y} = \text{class } k$$

$$\text{if } k = \arg \max_k q_k(x)$$

Soft max
regression

$$H(P, q)(x) = \left(-\sum_{k=1}^k P_k \log q_k \right) (x)$$

$$w_k^* = \arg \min_w E[H(P, q)(x)] \quad \left(\sum_{k=1}^k q_k = 1 \right)$$

$$w_{k+1} = w_{k+1} - \nabla_{w_k} E[H(P, q)(x)]$$

$m+1$ -dim.

$$\nabla_{w_k} E[H(P, q)(x)]$$

$$= E[(g_k - p_k)x]$$

$$\approx \frac{1}{n} \sum_i (g(x_i) - p(x_i)) x_i$$

$$\approx \frac{1}{n} \sum_i (g(x_i) - p(x_i)) x_i$$

$$E[X] q(x) = [0.1, 0.2, 0.3, 0.15, 0.25]$$

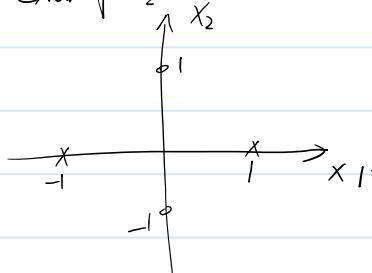
$$\Rightarrow \hat{y} = 3 \quad \text{3rd class}$$

5. Random forest and Boosting

1) CART (classification and regression trees)

④ classification tree., use a tree to represent a classifier

example:



classifier: $g = f(x)$.

input $x \in \mathbb{R}^m$ $m=2$.

output $g \in \{-1, +1\}$

$\begin{matrix} X \\ -1 \end{matrix}$

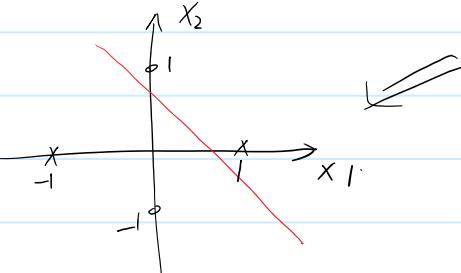
Past

$$g = \text{sign}(w^T x)$$

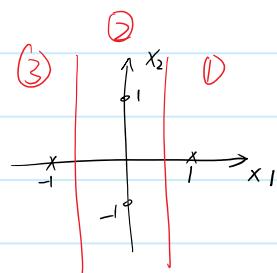
Now

Tree / Decision tree / classifier

$$f(\cdot)$$



Stump
fence.



(1)

$$y = -1$$

$$y = 1$$

$$y = -1$$

(3)

set indicator
function.

$$f(x) = \sum_i a_i I_{A_i}(x)$$

$$I_A(x) = \begin{cases} 0, & x \in A \\ 1, & x \notin A \end{cases}$$

$$= (-1)I_{A_1}(x) + (1)I_{A_2}(x) + (-1)I_{A_3}(x)$$

$$f(x) = 2x - 4$$

$$x_{n+1} = x_n - \epsilon(2x_n - 4)$$

$$= (1-2\epsilon)x_n + 4\epsilon$$

$$= (1-2\epsilon)^2 x_{n-1} + 4\epsilon + 4\epsilon(1-2\epsilon)$$

$$= (1-2\epsilon)^n x_0 + 4\epsilon(1+4\epsilon) + \dots + (1-2\epsilon)^{n-1}$$

$$= (1-2\epsilon)^n x_0 + 4\epsilon \cdot \frac{1-(1-2\epsilon)^n}{2\epsilon}$$

$$\therefore -1 < 1-2\epsilon < 1$$

when $n \rightarrow +\infty$ $(1-2\epsilon)^n$ or $(1-2\epsilon)^{n-1} \rightarrow 0$

$$\therefore x_{n+1} \approx 0 + 4\epsilon \cdot \frac{1}{2\epsilon}$$

$$= 2$$

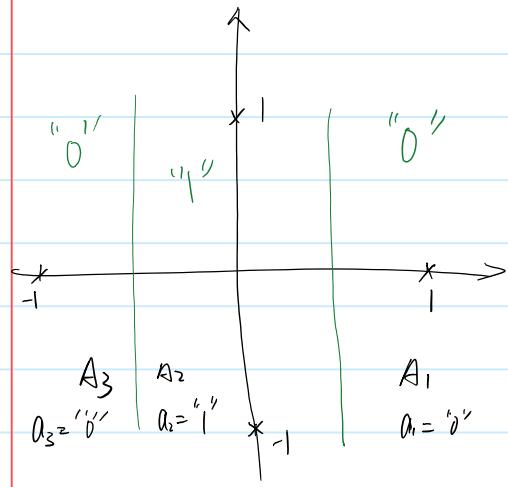
Week12

2017年11月20日 20:18

Review,

(highly nonlinear classification in CS way)

CART (classification & Regression tree)



partition of
input Space

each subset has a unique label

① classification tree.

- ① split of data
- ② along x_1

$$x_1 > \frac{1}{2}$$

$$\begin{array}{l} H = -\frac{1}{2} \log \frac{1}{2} \\ = \frac{1}{2} \end{array}$$

$$50\% "0" \quad 50\% "1"$$

100%

"0"

"1"

$$H = -1 \log 1 - 0 \log 0 = 0$$

Ave.

$$0 \cdot \frac{1}{4} + () \times \frac{3}{4}$$

$$\leq \frac{1}{4}$$

100%

"1"

"0"

$\frac{1}{3} "0"$

$\frac{2}{3} "1"$

$$H = -\frac{1}{3} \log \frac{1}{3} - \frac{2}{3} \log \frac{2}{3}$$

100% "0"

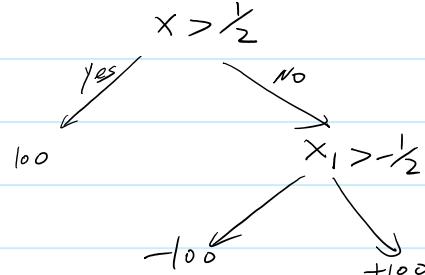
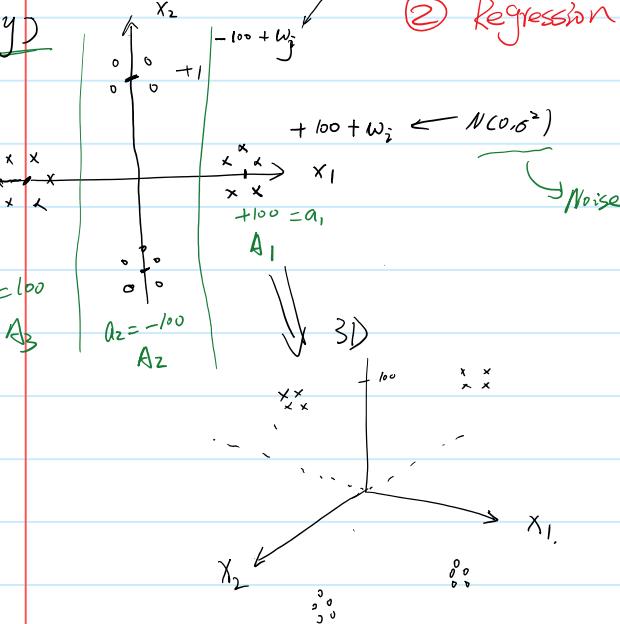
IDEA: reduce entropy

$$f(x) = \sum a_i I_{A_i}(x)$$

② Regression Tree.

regression function

$$y = f(x)$$



$$f(x) = \sum a_i I_{A_i}(x)$$

③ Training for CART.

$$D = \{(x_i, y_i)\} \implies f(\cdot)$$

dataset D $F(D)$

(1) classification, $F(D)$, "purity"

$\hat{p}_k^1 = \text{# of Point } (x_i, y_i)$

that belong to class k ($k=0, 1$)

$$F(D) = F(\{\hat{p}_k\})$$

$$= H(\{\hat{p}_k\})$$

$$= -\hat{p}_0 \log \hat{p}_0 - (1-\hat{p}_0) \log (1-\hat{p}_0)$$

(2) Regression

$$D = \{(x_1, y_1), \dots, (x_m, y_m)\}$$

$$f(x) = \frac{1}{m} \sum y_i$$

$$= \arg \min_a F(D)$$

$$F(D) = \min_a \sum_m (y_i - a)^2$$

\Rightarrow choose " a " to represent y .
a value

2) Random forest.

① Ensemble learning

$f_i()$: ^{independent} classifier i . $P_e < 50\%$

$$\text{majority vote classifier} = \text{majority } \left\{ \underbrace{\{f_i()\}}_{i=1} \right\}_{i=1}^{2n+1}$$

← odd number.

AS $n \rightarrow \infty$, Prob error $\rightarrow 0$

P_e (majority vote)

$$\text{Binomial distribution}, P_e = \sum_{k=n+1}^{2n+1} \binom{k}{2n+1} \varepsilon^k (1-\varepsilon)^{2n+1-k}$$

k out of $2n+1$
are in error.

$\varepsilon = P_e$ individual classifier. ex: 49%

② Apply to CART

$$D = \{(x_i, y_i)\} \rightarrow \text{random trees. } \{T_m\}$$

classification

$$T(X) = \text{Majority } (T_m(X))$$

How to generate
random trees

Regression

$$T(X) = \frac{1}{M} \sum_{m=1}^M T_m(X)$$

Random Trees: increase the randomness between individual classifier variables

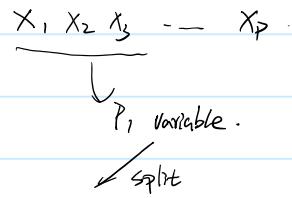
$T_m()$,

- Take a random subset of D

- when making a split decision

Only look at P_i random variables out of total P variables

along x_1 ,
or
 x_2



3) Boosting

① Basic idea,

$$P_c > 50\%$$

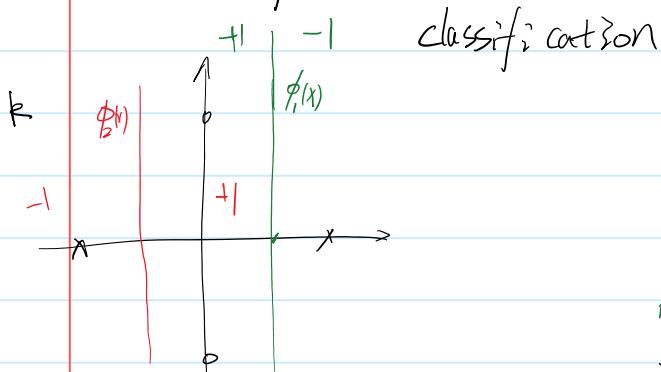
- a sequence of weak learners,

- combined in a strong learner

adding

$$P = 100\% \text{ (on training data)}$$

Example,



$\phi_1(x)$ (classifier 1)

$$= \begin{cases} -1, & x \geq \frac{1}{2} \\ +1, & x < \frac{1}{2} \end{cases}$$

$$P_e = \frac{1}{4} = 25\%$$

(try to previous mistake)

$\phi_2(x)$ (classifier 2).

$$= \begin{cases} -1, & x \leq -\frac{1}{2} \\ +1, & x > -\frac{1}{2} \end{cases}$$

$$P_e = \frac{1}{4} = 25\%$$

$$P_c = \frac{3}{4} = 75\% > 50\%$$

$f(x) = \phi_1(x) + \phi_2(x) = \begin{cases} 0, & x \geq \frac{1}{2} \\ 2, & -\frac{1}{2} \leq x \leq \frac{1}{2} \\ 0, & x \leq -\frac{1}{2} \end{cases}$

$$P_c = 100\%$$

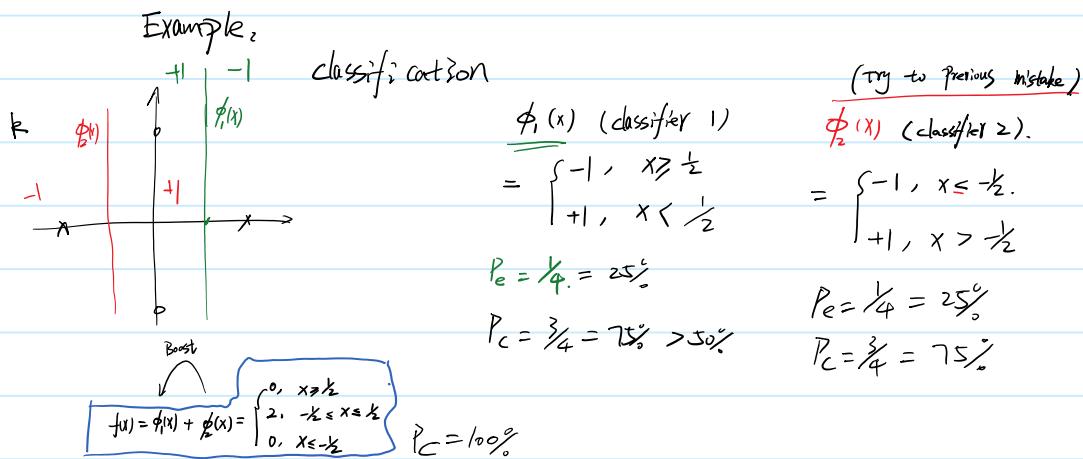
Review boosting

Basic idea.

 $f(x) = \text{classifier. / regression function}$

$$f(x) = \sum \alpha_m \phi_m(x)$$

$\hat{\alpha}$
Weak learner ($P_e < 50\%$)



(2) Ada Boost.

adaptive

 $f(x)$: classification function. $f_m(x) \longrightarrow f(x)$

$$\phi_m(x, y_m) = \begin{cases} -1, & x > \frac{1}{2} \\ +1, & x < \frac{1}{2} \end{cases}$$

$$y_m = [\frac{1}{2}, 0]^T$$

$$\alpha_m = 1$$

$$f_m(x) = f_{m-1}(x) + (\alpha_m \phi_m(x, y_m))$$

How to find α_m, γ_m .

Given, $f(x) - \phi_m(x)$

$$(\alpha_m, \gamma_m) = \underset{\alpha, \gamma}{\operatorname{argmin}} \mathcal{L}(y, f_m(x))$$

\mathcal{L} → loss function.

$y = f(x)$ = out of classifier (truth) $\in \{-1, +1\}$

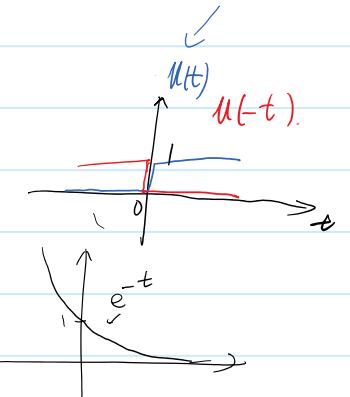
$$f_m(x) = \text{"boost" classifier} = f_{m-1}(x) + \alpha \phi_m(x, \gamma) \in \{-1, +1\}$$

not derivative

$\mathcal{L}(y, f_m(x))$ = loss function

$$\text{ideally} = \begin{cases} 0 & \cdot y = f_m(x) \\ 1 & \cdot \text{o.w.} \end{cases} = \mathcal{L}(-y f_m(x))$$

$$\approx e^{-y f_m(x)}$$



⑧ AdaBoost Optimization.

$$D = \{(x_i, y_i)\}$$

$$E[\mathcal{L}(y, f_m(x))] = \mathcal{L}_m = C \cdot \sum_{i=1}^n \mathcal{L}(y_i, f_m(x_i))$$

$$= C \cdot \sum e^{-y_i f_m(x_i)}$$

$$= C \cdot \sum e^{-y_i f_{m-1}(x_i) - y_i \alpha \phi(x_i, \gamma)}$$

See textbook

$$\underset{\alpha, \gamma}{\operatorname{argmin}} \mathcal{L}_m \rightarrow \alpha_m, \gamma_m$$

basic idea, optimize α first
and optimize γ 2nd.

④ Gradient Boosting

see in Elements of statistical learning

In classification / regression.

$$D_{\text{training data}} = \{(x_i, y_i)\}$$

→ $f(x)$

y_i output / ground truth.

Want to find $f(x)$ st.

$E[\mathcal{L}(y, f(x))]$ is minimized.

$$E[\mathcal{L}(y, f(x))] = \frac{1}{n} \sum \mathcal{L}(y_i, f(x_i))$$

function

$$f^*(x) = \underset{f(x)}{\arg \min} E[\mathcal{L}(y, f(x))]$$

<u>function</u>	<u>Variate</u>	X
<u>approximate</u>		✓

$$\vec{f} = [f(x_1), f(x_2), \dots, f(x_n)]^T \quad \underline{E[\mathcal{L}(y, f(x))]}$$

not known
but know what it
should be

$$\begin{aligned} \vec{y} &= [y_1, y_2, \dots, y_n]^T \\ &\approx \frac{1}{n} \sum \mathcal{L}(y_i, f(x_i)) \\ &= \frac{1}{n} \sum \mathcal{L}(y_i, f_i) \\ &= \mathcal{L}(\vec{y}, \vec{f}) \end{aligned}$$

find \vec{f} by GD

$$\vec{f}_m = \vec{f}_{m-1} - \sum_m \nabla_{\vec{f}} \mathcal{L} |_{\vec{f} = \vec{f}_{m-1}}$$

find \vec{f} by GD

$$\vec{f}_m = \vec{f}_{m-1} - \sum_m \nabla_{\vec{f}} \mathcal{L} |_{\vec{f}=\vec{f}_{m-1}}$$

$$\vec{f}_m \xrightarrow[m \rightarrow \infty]{} \vec{f}^* = \underbrace{[\vec{f}^*(x_1) - \vec{f}(x_1)]^\top}_{\text{optimizing values of function}}$$

optimizing values of function

not function.

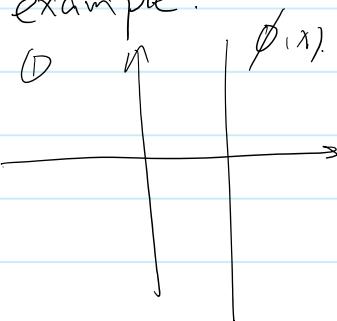
new $x \rightarrow f^*(x) = ?$
 $\notin \{x_1, x_2, \dots, x_n\}$

extra process $\xrightarrow{\text{introducing function.}}$ interpreting the gradient function

$$\phi(x) = \phi(x; \gamma)$$

\uparrow parameters

example.



ϕ, f i

not need to
be Gaussian
can be any
function.

$$\textcircled{2} \quad \phi(x; \gamma) = c \cdot e^{-\frac{(x-m)^2}{2\sigma^2}}, [m, \sigma^2]^\top = \gamma$$

minimize gradient g & $\phi(x; \gamma)$
at points x_i

$$\gamma_m = \arg \min_{\gamma} \sum_i \mathcal{L}(g_i, \phi(x_i; \gamma))$$

$$g_i = \nabla_{\vec{f}} \mathcal{L}(\vec{g}, \vec{f}) (\vec{f} = \vec{f}_{m-1})$$

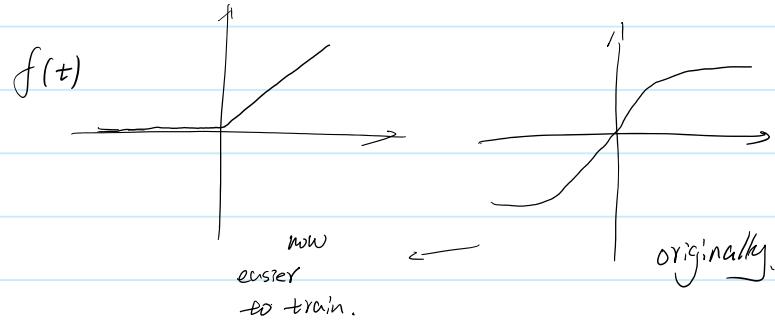
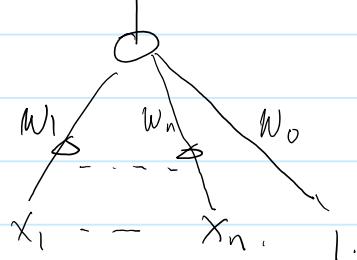
$$\Rightarrow f_m(x) = f_{m-1}(x) - \varepsilon \phi(x, \gamma_m)$$

example,

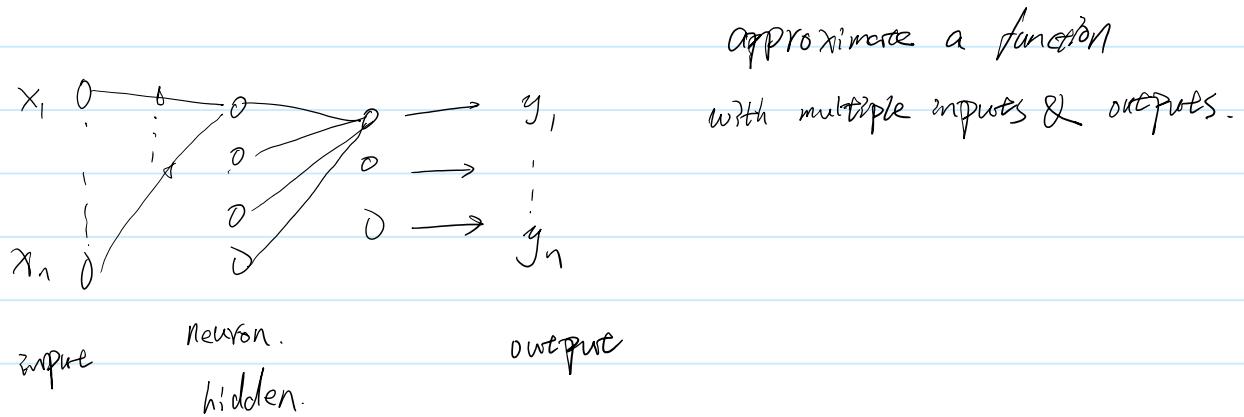
$$\mathcal{L}(g, f(x)) = (g - f(x))^2 \Rightarrow L_2 \text{ boost}$$

5 Neural Networks

1) a "neuron" $y = f(w_0 + w_1 x_1 + \dots + w_n x_n)$
 $= f(w^T x)$



2) A neural network



3) Training

$$D = \{(x_i, y_i)\} \rightarrow \text{a neural network (fixed structure)}$$

\Downarrow
 $f(x)$

$$\text{neural network} = f(x) = f(x_i - w)$$

training $w_t = w_{t-1} - \epsilon \nabla_w \sum L(y_i, f(x_i))$ | $w = w_{t-1}$

Adam algorithm
 \Rightarrow backpropagation

kingma

Adam

algorithm

Problem.

(1) converge is too slow

(2) converge to a plane ?

pseudo-newton technique // not memory efficient

Week2

2018年1月29日 18:00

Review

Machine Learning

- classification $x \xrightarrow{f(\cdot)} \text{label / class}$
- Regression $x \xrightarrow{f(\cdot)} y \quad P[x \in \text{class } k] = f(x)$

Gaussian Distribution.

$$X \sim N(m, \Sigma) \quad P(X) = (2\pi)^{-\frac{n}{2}} |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2} \|X-m\|^2 / \Sigma}$$

Marginal & Conditional

Marginal, $\begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} \quad X_i, \quad X_i \sim N(m_i, \Sigma_{ii})$

$$\Sigma = \begin{bmatrix} m_1 & & \\ m_2 & \Sigma_{11} & \Sigma_{12} \\ \vdots & & \\ m_n & \Sigma_{21} & \Sigma_{22} \end{bmatrix}$$

Conditional

$$X_1 | X_2 \sim N(m_{12}, \Sigma_{12})$$

$$m_{12} = m_1 + \Sigma_{12} \Sigma_{22}^{-1} (X_2 - m_2)$$

$$\Sigma_{12} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$P(X_1 | X_2) = \frac{P(X)}{P(X_2)} \leftarrow P(X_1, X_2)$$

$$= \frac{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2} \|X-m\|^2 / \Sigma}}{(2\pi)^{\frac{n}{2}} |\Sigma_2|^{\frac{1}{2}} e^{-\frac{1}{2} \|X_2-m_2\|^2 / \Sigma_{22}}}$$

$$\rightarrow \|X_1 - m_{12}\|^2 / \Sigma_{12}$$

$$f(u, v) = u^2 + auv + v^2 + bu + cv + d. \quad \frac{\partial f}{\partial u} = 2u + av + b = 0 \quad u = \frac{-av-b}{2}$$

$$f(u, v) = (u - m(v))^2 + \tilde{\Sigma}(v) \quad \frac{\partial f}{\partial u} = 2(u - m) = 0 \Rightarrow u = m$$

4) Optimal classification of Gaussian Distribution

(1) Two class problem. (hard classifier)

$$H_0: \mathbf{x} \sim N(\mathbf{m}_0, \Sigma_0)$$

$$H_1: \mathbf{x} \sim N(\mathbf{m}, \Sigma_1)$$

optimal classifier

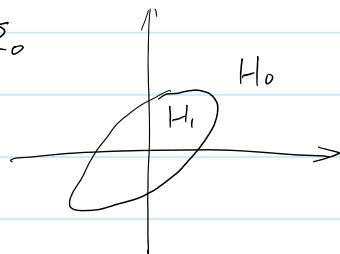
$$\log \frac{P(x|H_1)}{P(x|H_0)} \stackrel{H_1}{\geqslant} \log \frac{P_0}{P_1} \quad \text{L-LR}$$

Log likelihood ratio

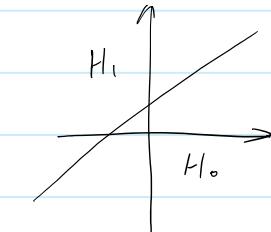
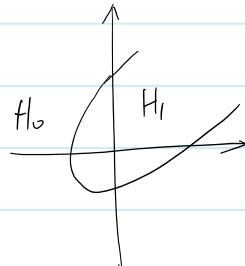
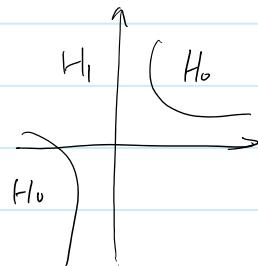
$$\begin{aligned} LHS &= \log \frac{(2\pi)^{-\frac{D}{2}} |\Sigma_1|^{-\frac{1}{2}} e^{-\frac{1}{2} \|x - m_1\|_{\Sigma_1}^2}}{(2\pi)^{-\frac{D}{2}} |\Sigma_0|^{-\frac{1}{2}} e^{-\frac{1}{2} \|x - m_0\|_{\Sigma_0}^2}} \\ &= C - \frac{1}{2} (\|x - m_1\|_{\Sigma_1}^2 - \|x - m_0\|_{\Sigma_0}^2) \\ &= \text{quadratic in } x \end{aligned}$$

Boundary x , $LHS = 0$.Shape
Depends on

$$\Sigma_1, \Sigma_0$$



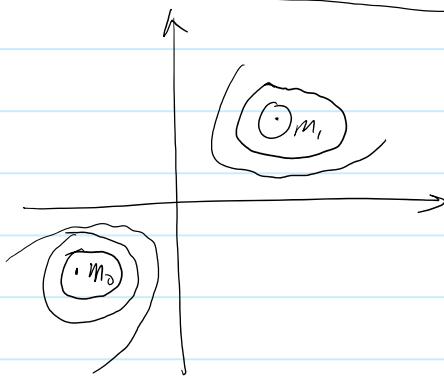
$$\frac{\|x - m_1\|_{\Sigma_1}^2 - \|x - m_0\|_{\Sigma_0}^2 + C'}{2} = 0$$

Quadratic in x .

Special case.

$$\textcircled{1} P_0 = P_1 = \frac{1}{2}$$

$$\Sigma_0 = \Sigma_1 = \sigma^2 I$$

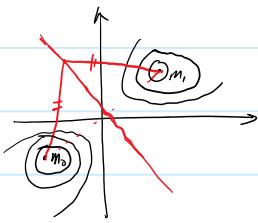


$$\|x - m\|_{\Sigma}^2 = \frac{1}{\sigma^2} \|x - m\|^2$$

$$\|x - m_1\|^2 \stackrel{H_0}{\gtrless} \|x - m_0\|^2.$$

Note, '-' sign in front.

$$\|x - m_1\| \stackrel{H_0}{\gtrless} \|x - m_0\| \Rightarrow \text{min distance classifier}$$



$$\textcircled{2} \quad P_0 = P_1, \quad \Sigma_0 = \Sigma_1 = \sigma^2 I \quad \Rightarrow \quad m_1^T x \stackrel{H_1}{\gtrless} m_0^T x \quad \|x - m_1\|^2 = \|x\|^2 - 2m_1^T x + \|m_1\|^2.$$

$$\|m_0\|^2 = \|m_1\|^2$$

& correlation classifier

② "soft" classifier

Two-class

$$H_0: x \sim N(m_0, \Sigma_0) \rightarrow P(x | H_0)$$

$$H_1: x \sim N(m_1, \Sigma_1) \rightarrow P(x | H_1)$$

$$\text{"soft" classifier} \Rightarrow P[H_0 | x] = \frac{P(x | H_0) P[H_0]}{P(x)}$$

$$\downarrow \quad \quad \quad \downarrow \quad \quad \quad P_0$$

$$P(x | H_0) P[H_0] + P(x | H_1) P[H_1]$$

$$P[H_0 | x] = \frac{1}{1 + \frac{P(x | H_1) P_1}{P(x | H_0) P_0}}$$

$$\leftarrow \frac{\frac{1}{|\Sigma_1|^{\frac{N}{2}}} e^{-\frac{1}{2}(\|x - m_1\|_{\Sigma_1}^2 - \|x - m_0\|^2)}}{|\Sigma_0|^{\frac{N}{2}}}$$

Quadratic in x .

$$\|x - m\|^2 = \|x\|^2 - 2m^T x + \|m\|^2 \underset{m_0}{\leftarrow}$$

Special case.

$$P_0 = P_1$$

$$\Sigma_0 = \Sigma_1 = \sigma^2 I$$

$$P[H_0|x] = \frac{1}{1 + e^{-\frac{1}{2\sigma^2}((2(m_0 - m_1)^T x + \|m_0\|^2 + \|m_1\|^2)/2)}} \\ = \frac{1}{1 + e^{-w^T x}} \left[\begin{array}{c|c} x & 1 \end{array} \right] \text{ augmented.}$$

Boundary, for general case.
Quadratic

$$\Rightarrow \text{logistic regression } g = f(x) = P[x \in H_0] = P[H_0|x]$$

$$\text{Boundary}, P[H_0|x] = \frac{1}{2}.$$

$$\Rightarrow w^T x = 0 \text{, linear equation}$$

5) Optimal Regression w/ Gaussian

① Optimal Regression

$$\text{Input } x \rightarrow N(m_x, \Sigma_{xx})$$

$$\text{Output } y \rightarrow N(m_y, \Sigma_{yy})$$

$$(x, y) \sim N\left(\begin{bmatrix} m_x \\ m_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}\right)$$

Regression Problem.

Estimate y from x .

$$\text{Optimal estimate} = E[y|x]$$

Data

$$x_i \rightarrow \hat{y}_i = E[y|x_i] = f(x_i)$$

$$= m_y + \Sigma_{yx} \Sigma_{xx}^{-1} (x - m_x)$$

$$= \Sigma_{yx} \Sigma_{xx}^{-1} x + (m_y - \Sigma_{yx} \Sigma_{xx}^{-1} m_x)$$

$$= w^T x \leftarrow \text{Augmented } \left[\begin{array}{c|c} x & 1 \end{array} \right]$$

\Rightarrow Linear regression

$w^T \Rightarrow \text{matrix } x$
when y has multi-D

$$\underset{\text{est. error}}{\text{definition: }} \mathbb{E} \left[(y - \mathbb{E}[y|x]) (y - \mathbb{E}[y|x])^\top \right]$$

$$\Sigma_{y|x} = \Sigma_{yy} - \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$$

\Rightarrow How good the estimation is?

error covariance matrix

MSE covariance matrix.

b) "Training" for Gaussian models:

① Classification

$$H_0: X \sim N(m_0, \Sigma_0)$$

$$H_1: X \sim N(m_1, \Sigma_1)$$

training, parameters estimation



$$D \rightarrow m, \Sigma$$

② Regression,

$$D = \{x_1, x_2, \dots, x_n\}$$

$$X \sim N(m_x, \Sigma_{xx})$$

$$y \sim N(m_y, \Sigma_{yy})$$

$$(x, y) \sim N \begin{bmatrix} m_x \\ m_y \end{bmatrix}, \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix}$$

MLE.

$$(\hat{m}, \hat{\Sigma}) = \underset{m, \Sigma}{\operatorname{argmax}} \log P(D|m, \Sigma)$$

Gaussian

$$\hat{m} = \frac{1}{n} \sum x_i = \text{Sample mean}$$

$$\hat{\Sigma} = \frac{1}{n} \sum (x_i - \hat{m})(x_i - \hat{m})^T$$

$$P(m, \Sigma) \cdot \text{Prior}$$

MAP: (max a posteriori) estimate

$$\hat{m}, \hat{\Sigma} = \underset{m, \Sigma}{\operatorname{argmax}} \left\{ \log P(D|m, \Sigma) \right\}$$

Guassian
without distribution

$$= \operatorname{argmax} \{ \log P(D|m, \Sigma) + \log P(m, \Sigma) \}$$

Useful. When the amount of data

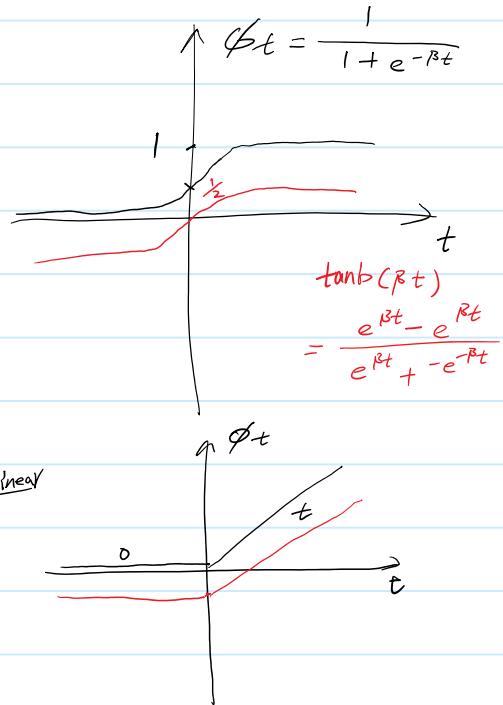
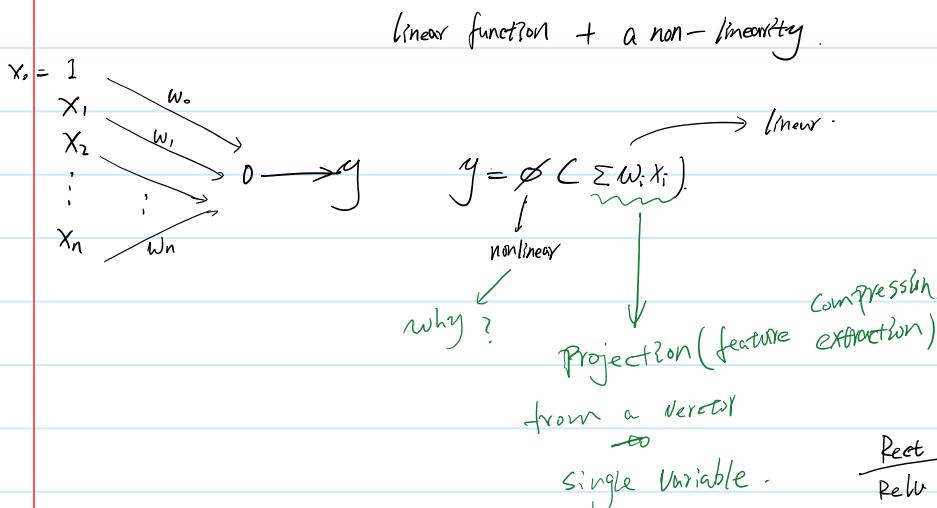
is not large.

clustering

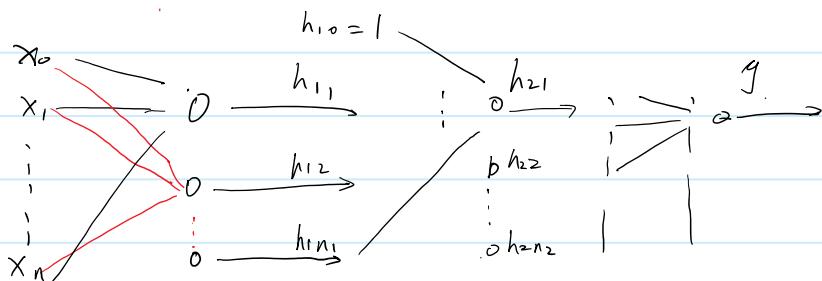
3. Neural Networks

1) Neural networks

① A "neuron"



② A neural network.



$$x \xrightarrow{w_1} h_1 \xrightarrow{w_2} h_2 \dots \xrightarrow{w_n} y$$

$$\mathbb{R}^n \quad \mathbb{R}^{n_1} \quad \mathbb{R}^{n_2} \quad \dots \quad \mathbb{R}$$

Set up, fully connected

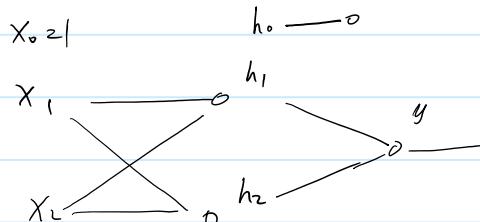
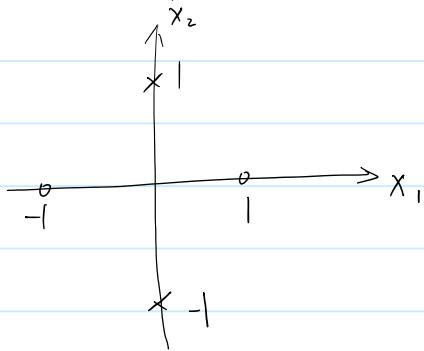
feed forward

N.N
F.F
RNN

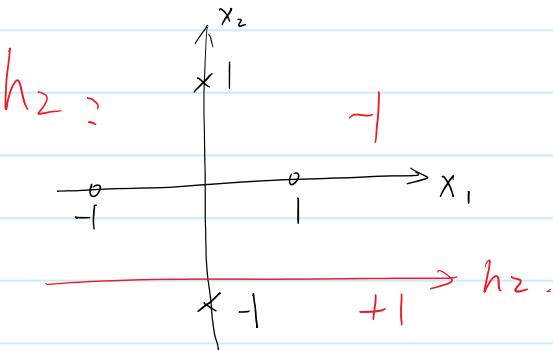
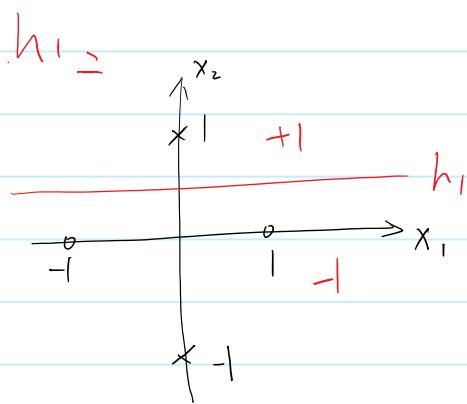
Recurrent
(has feedback)

MLP
CNN
AE

(3) XOR Problem.



$$\mathbb{R}^2 \rightarrow \mathbb{R}^2$$



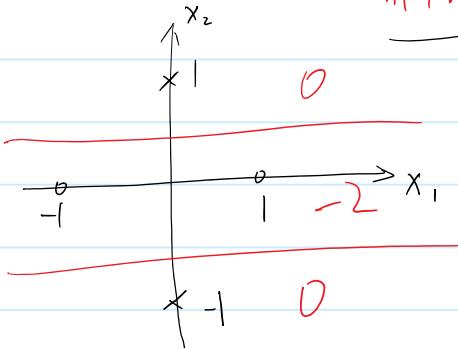
Non-linear function

$$h_1(x_1, x_2) = \begin{cases} +1 & x_2 > \frac{1}{2} \\ -1 & x_2 < \frac{1}{2} \end{cases}$$

$$= \phi(w_1 x_1 + w_2 x_2 + w_0) \quad \leftarrow \begin{array}{c} \phi(t) \\ + \\ t \end{array}$$

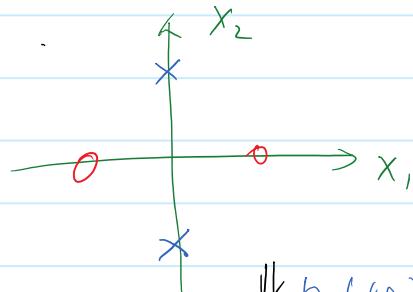
before
 $h_1 + h_2$

$h_1 + h_2$

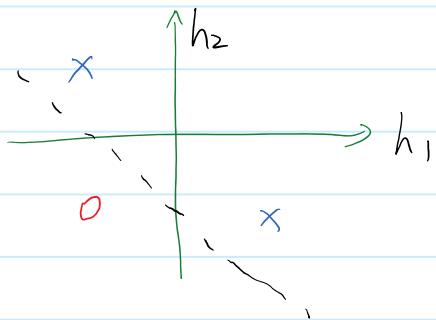


$$h_1 + h_2 = y$$

linear



after



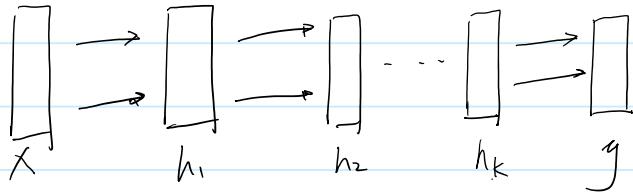
Activation function

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU) [2]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) [3]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

Week3

2018年2月5日 17:08

Review, Neural networks



$$h_k = \phi(w^{k-1} h_{k-1})$$

↗ non-linear
 ↓ matrix
 ↘ column vector
 componentwise

$$\phi\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_n) \end{pmatrix}$$

2) Function Approximation by NN.

① function approx. est. problem.

function, $g = f(x)$

$f(\cdot)$ is unknown

task, $D = \{(x_i, y_i)\}$

Function est. / appr. $\hat{f}(\cdot)$.

$D \rightarrow \hat{f}(\cdot)$

interpolation
regression

①

Example. LR: $\hat{f}(x) = w^T x$ $\xrightarrow{\text{linear function}}$ n^{dim} # of terms $O(n)$

Adv.
Simple
Scales well w/ dim
Prob.
large error
when $f(\cdot)$ is nonlinear.

$$w = E[x x^T]^{-1} E[x y]$$

$$\approx \left(\frac{\sum x_i x_i^T}{n} \right)^{-1} \approx \frac{\sum x_i y_i}{n}$$

(2) Example: polynomial problem.

$$(D = f(x) = a_0 x^k + a_1 x^{k-1} + \dots + a_k)$$

Adv: simple; nonlinear for 1D problem

$$(a_0^*, a_1^*, \dots, a_k^*) = \arg \min_{a_0, \dots, a_k} E[(y - f(x))^2]$$

Problem: doesn't scale well for n-dim.

2D polynomial of order k ($k=2$)

$$f(x) = ax_1^2 + bx_2^2 + cx_1 + dx_2 + ex_1x_2 + f$$

$$f(x) = ax_1^2 + bx_2^2 + \underbrace{cx_1}_{\alpha_0=0, \alpha_1=1, \alpha_2=0} + \underbrace{dx_2}_{\alpha_0=1, \alpha_1=0, \alpha_2=0} + ex_1x_2 + f$$

$\alpha_0 = 0$
 $\alpha_1 = 1$
 $\alpha_2 = 0$
 $\alpha_0 = 1$
 $\alpha_1 = 0$
 $\alpha_2 = 1$
 $\alpha_0 = 0$
 $\alpha_1 = 2$
 $\alpha_2 = 0$
 $\alpha_0 = 0$
 $\alpha_1 = 0$
 $\alpha_2 = 2$

n-dim polynomial order k.

$$f(x) = \sum x_1^{\alpha_1} x_2^{\alpha_2} \dots x_n^{\alpha_n} + c$$

$$\alpha_0 + \alpha_1 + \alpha_2 + \dots + \alpha_n = k.$$

$$\alpha_0 \geq 0, \alpha_1 \geq 0, \alpha_2 \geq 0 \dots \alpha_n \geq 0$$

Stars & Bars

of terms $O(n^k)$

$$n = 256 \times 256 \rightarrow (2^{16})^{10} = 2^{160}$$

(3) Example: Fourier transform.

$$1D: f(x) = \sum_{m=0}^{M-1} f_m e^{j w_0 m x} \quad \# \text{ of terms: } O(M)$$

$$2D: f(x) = \sum_{m_1, m_2} f_{m_1, m_2} e^{j w_0 (m_1 x_1 + m_2 x_2)} \quad \# \text{ of terms: } O(M^2)$$

$$nD: \# \text{ of terms } O(M^n)$$

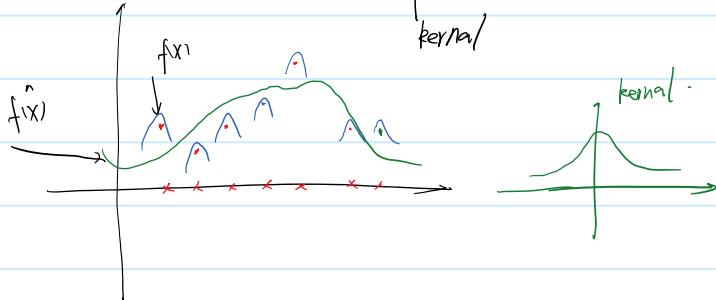
Adv: Fourier can represent any formulation

Problem: only in low dim problems

④ Example: kernel regression (sum)

$$f(x) = \sum w_i \phi(x - x_i)$$

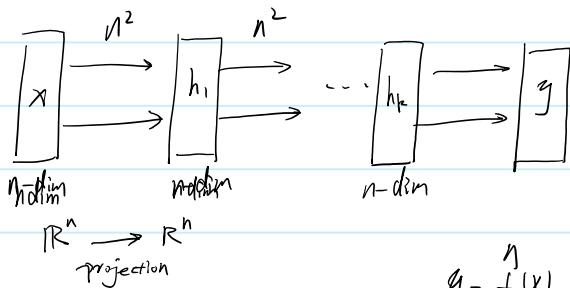
need training data every time.
 x_i



n -dim, 100^n point.

$n = 100 \rightarrow 100^{100}$ point.

⑤ Example: NN



$\begin{matrix} \text{n-dim column vector} \\ h_i = \phi(w^T x) \\ \uparrow \\ \text{n-dim column vector} \end{matrix}$
 w $\in \mathbb{R}^m$
 $m \times n$ matrix
 $O(n^2)$

Total # of weights,

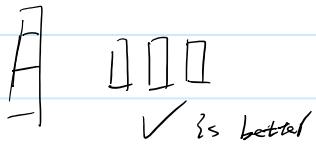
$$\begin{aligned} g &= f(x) \\ &= w^T h_k(x) \end{aligned}$$

(orthogonal)

NN's ability to approximate

→ only one hidden layer is sufficient
 (may need an exponential # of nodes)

For some nodes,



✓ is better

n input, n hidden layer nodes



$$(n/3 \times n) + (n/3)^2 + (n/3)^2$$

$$= \frac{n^2}{3} + \frac{n^2}{9} + \frac{n^2}{9} = \frac{5}{9}n^2$$

3) Neural Network training

① Problem:

$$D \rightarrow w$$

② Optimization

$$w^* = \arg \min_w J(w)$$

example, MSE

$$\begin{aligned} J(w) &= E[(g - f(x))^2] \\ &\quad \text{out put} \quad \text{NN} \quad \text{input to NN} \\ &= E[(g - f(x/w))^2] \end{aligned}$$

$$= \frac{1}{N} \sum_{i=1}^N (y_i - f(x_i/w))^2.$$

example, cross entropy

random variable X , X discrete

$$\begin{array}{c|cccc} X & x_1, x_2, \dots, x_n \\ \hline P & p_1, p_2, \dots, p_n \end{array} \Rightarrow \text{true PD}$$

Probability Density

approximation to P , q

$$\begin{array}{c|cccc} X & x_1, \dots, x_n \\ \hline q & q(x_1), \dots, q(x_n) \end{array}$$

$$H(P, q)$$

$$= E_p[-\log q_{x_i}]$$

$$= -\sum p_k \log q(x_k)$$

$$= -\sum p_k \log q_k$$

$$\text{in NN, } q = q(x) = f(x/w)$$

$$H(P, q) \approx - (q_i \log f(x_i/w) + (1-q_i) \log (1-f(x_i/w)))$$

binary case

$$\Leftrightarrow H(P, q) = H(P) + D(P // q)$$

entropy
of
 X

k -L distance
from P to q .

$D(P // q) \geq 0$
 $= 0$ when $P = q$

$$-\sum p_k \log p_k \quad \sum p_k \log \frac{p_k}{q_k}$$

Example: likelihood function,

$$\omega = \underset{\omega}{\operatorname{argmin}} \mathbb{E} [-\log g(y|x, \omega)]$$

g : conditional pdf of y given x .

ω : parameters

$$\mathbb{E} [-\log g(y|x, \omega)]$$

$$\approx \frac{1}{N} \sum_{i=1}^N (-\log g(y_i|x_i, \omega))$$

Example: $g(y|x, \omega)$

$$y|x, \omega \sim N(m, \sigma^2)$$

$$m = \hat{f}(x|\omega)$$

\nwarrow
Neural network

$$\Rightarrow \mathbb{E} [-\log g(y|x, \omega)]$$

$$= C + \mathbb{E} \left[\frac{(y - \hat{f}(x|\omega))^2}{2\sigma^2} \right]$$

= MSE

Example: same as before.

$$m = \hat{f}(x|\omega) = \text{a NN}$$

$$\sigma^2 = \hat{g}(x|\omega) = \text{another NN}$$

$$\mathbb{E} [-\log g(y|x, \omega, \omega')] = C + \mathbb{E} [\log g(x|\omega')] + \mathbb{E} \left(\frac{(y - \hat{f}(x|\omega))^2}{2\hat{g}(x|\omega')} \right)$$

$$\approx C + \frac{1}{N} \sum \left(\log \hat{g}(x|\omega') + \frac{(y_i - \hat{f}(x_i|\omega))^2}{2\hat{g}(x_i|\omega')} \right)$$

→ Criterion

→ optimization.

② Optimization.

→ traditional opt. technique (e.g. gradient descent)

→ Problems specific to deep structure (nested structure)

→ fitting the training data. vs. generalization.

overfit ⇒ regularization
⇒ robust estimation

Week 5

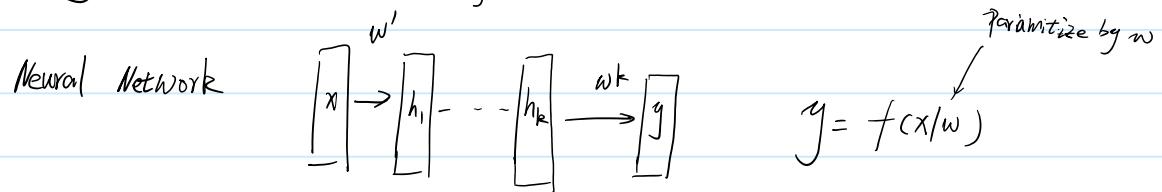
2018年2月19日 17:09

3) (Deep Neural network training)

- ① General problem
- ② Criteria
- ④ Review of optimization
- ⑤ Issues related to Deep neural network.

① General problem

Training data $D = \{(x_i, y_i)\}$



Weight $[w', w^2, \dots, w^k] = w$

Training $D \rightarrow w$

Solved as a optimization problem.

$$w^* = \underset{w}{\operatorname{argmin}} J(D/w)$$

② Examples of $J(D/w)$

$$\overbrace{\quad}^{E[y(x, y)/w]}$$

example, MSE

$$J(x, y/w) \rightarrow E[(y - f(x/w))^2]$$

$$\approx c \sum (y_i - f(x_i/w))^2$$

Example: Cross Entropy

$$\begin{aligned}
 J(x, y | w) &= E[-\log f(x|w)] \\
 &\stackrel{\text{true } P[y=1]}{\downarrow} \quad \stackrel{\text{neural network output}}{=} P[y=1] \text{ (estimated)} \\
 &= -P(x) \log f(x|w) - (1-P(x)) \log (1-f(x|w)) \\
 &\approx c \sum [-y_i \log f(x_i|w) - (1-y_i) \log (1-f(x_i|w))]
 \end{aligned}$$

Example: (conditional) likelihood function

Data: $P(x)$, pdf of n x

$$P(x) = P(x|\theta) \quad \underbrace{\text{parameters}}$$

likelihood function.

$D = \{x_1, x_2, \dots, x_n\}$ data and samples from X .

likelihood function, $\tilde{P}(D|\theta) = \tilde{P}(x_1, x_2, \dots, x_n | \theta)$

Example: $P(x|\theta)$ = Gaussian w/ (m, σ^2)

$$\theta = (m, \sigma^2)$$

$$P(\theta | D) = P(x_1, \dots, x_n | \theta)$$

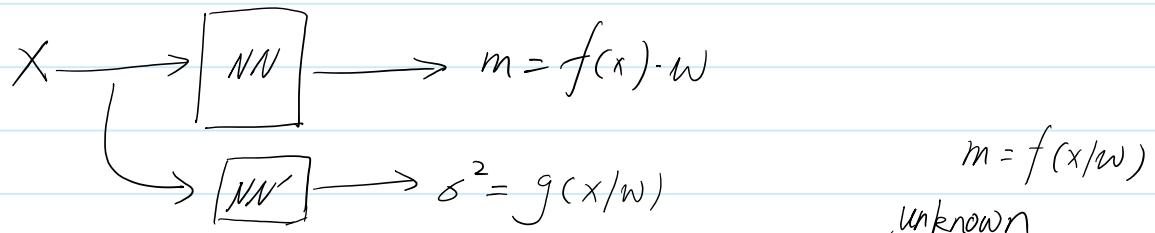
$$\pi(P(x_i | \theta)) = \pi \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - m)^2}{2\sigma^2}}$$

$$\text{MLE}, \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log \tilde{P}(D|\theta)$$

Gaussian example,

$$\hat{m} = \frac{1}{n} \sum x_i$$

Applied to Neural Network



likelihood function $P(y|x, \theta)$ = Gaussian (m, s^2)

any density function $= P(y|x, m) = P(y|x, f(x/w))$

NN training as MLE

likelihood function $P(y_1, \dots, y_n | x_1, \dots, x_n, m)$

$$= \prod_i \frac{1}{\sqrt{2\pi s^2}} e^{-\frac{(y_i - f(x_i/w))^2}{2s^2}}$$

$\text{MLE}, w^* = \arg \max_w P(y_1, y_2, \dots, y_n | x_1, \dots, x_n, m)$
 $= \text{non-linear function of weights.}$

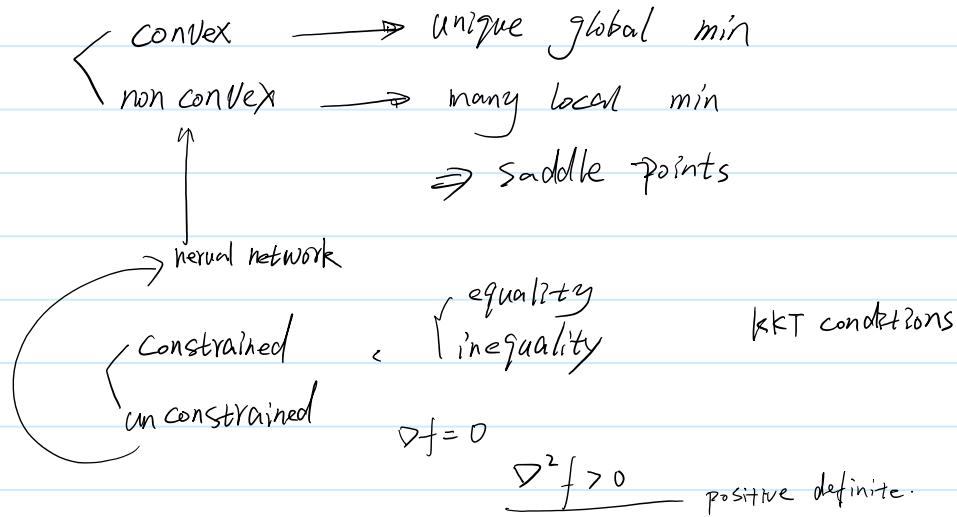
(Same as MSE)

$\text{MLE} \Rightarrow$ ① enlarge the training criteria.

ex, s^2 unknown \rightarrow know

② generate powerful density function

4) Review. optimization methods.



① Simplest case

$$x^* = \arg \min_x f(x)$$

necessary condition
not sufficient

$$\nabla f(x) = 0$$

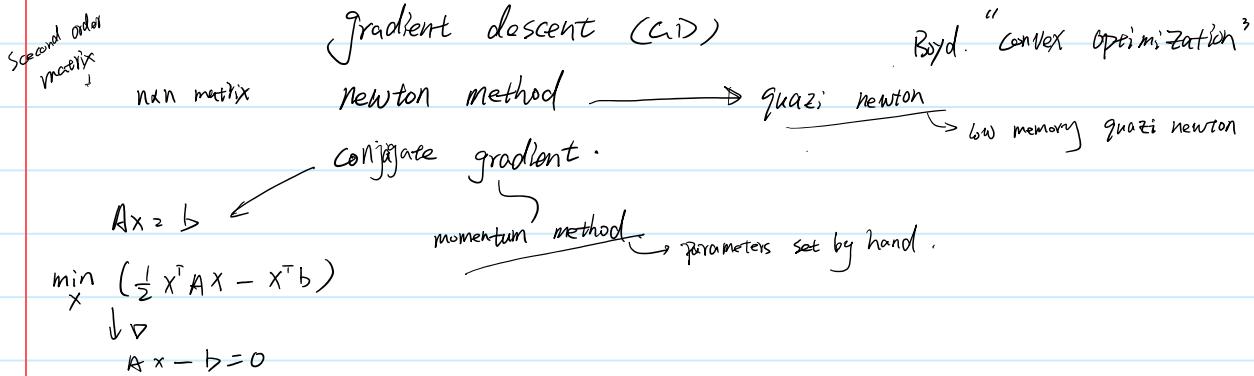
$$\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

Problem: $\nabla f(x)$ often can not be solved analytically.

Solution: iterative solution

"

② Iterative solutions:



③ newton method

$$x_0 \rightarrow x_1 \rightarrow x_2 \cdots x^*$$

(recall $x^* = \min_x f(x)$)

$$x f(x) = 0$$

$$g(x) = \nabla f(x)$$

Given x_n

Find x_{n+1}

$$g(x_{n+1}) \quad \text{ideally} \quad g(x_{n+1}) = 0$$

$$\text{Taylor expansion} \quad \frac{\partial T}{\partial x} \quad g(x_{n+1}) \approx g(x_n) + \left(\frac{\partial g(x_n)}{\partial x} \right) (x_{n+1} - x_n)$$

↓
 Vector

↗
 matrix

$$g(x_n) + \nabla g(x_n)(x_{n+1} - x_n) = 0$$

$$x_{n+1} - x_n = -(\triangledown g(x_n))^{-1} g(x_n)$$

$$= - \left(D^2 f(x) \right)^{-1} \left(Df(x_n) \right)$$

Hessian matrix

$$\text{Actual algorithm.} \quad \downarrow^{\text{line search}}$$

$$x_{n+1} = x_n - \alpha \left(D^2 f(x_n) \right)^{-1} (Df(x_n))$$

$$\nabla^2 f(x_n) = \begin{bmatrix} -\frac{\partial^2 F}{\partial x_i^2} & \frac{\partial F}{\partial x_i x_i} \\ \vdots & \ddots \\ \frac{\partial F}{\partial x_n x_n} & \frac{\partial F}{\partial x_n x_j} \end{bmatrix}$$

Problems : dim is high

Solutions \rightarrow quasi, low-mem Newton method.

→ drastic approximation

$$(\nabla^2 f)^{-1} \approx I$$

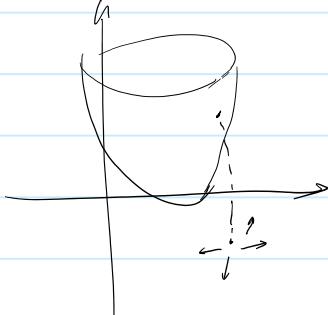
$$\Rightarrow x_{n+1} = x_n - \alpha Df(x_n) \Rightarrow GD$$

④ Gradient Descent

$$x_{n+1} = x_n - \alpha \nabla f(x_n)$$

Learning rate
step size.

direction fast descent.
 $\rightarrow f(x_n)$

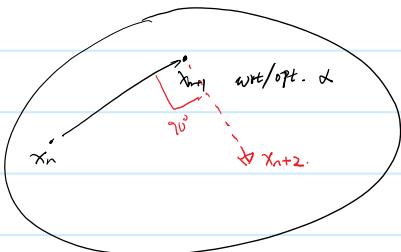


① fast change $\nabla f(x)$
 $\rightarrow f(x)$

② $\rightarrow f(x)$ descent.

line search,

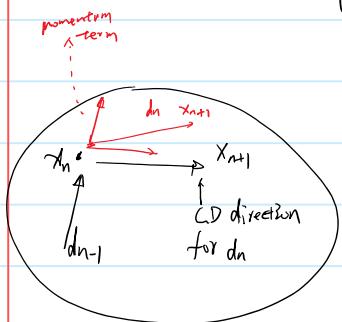
Select α , $f(x_{n+1}) = f(x_{n+1}(\alpha))$ is min wrt α .



⑤ Conjugate gradients method.

Basic idea.

d_{n-1} direction generated x_n .



$$x_n = x_{n-1} + \alpha d_{n-1}$$

$$d_n = -\nabla f(x_n) + \beta_n d_{n-1}$$

$$CD: d_{n-1} = -\nabla f(x_{n-1})$$

orthogonal

if $v^\top v = 0$ then $v^\top v = 0$, $v^\top v = 0$

conjugate $v^\top A v = 0$

A-orthogonal

A, SPD matrix

Symmetric Positive

$$d_n^T \nabla^2 f(x_n) d_{n-1} = 0 \quad \text{Conjugate relation}$$

$$\beta_n = \frac{(\nabla f(x_n) - \nabla f(x_{n-1}))^T \nabla f(x_{n-1})}{\nabla f(x_{n-1})^T \nabla f(x_{n-1})}$$

G. Strang, Applied math.
? , numerical Recipe. in C

optimal point x^*

$$\Delta d_k A d_{k-1} = 0 \quad \dim(x) = n.$$

$$\underbrace{d_i^T A d_j = 0}_{i \neq j} \quad n \text{ d}_k \text{'s. ; linearly independent.}$$

$$x^* = \sum x_i d_i \quad \Rightarrow \quad d_1 \rightarrow d_2 \rightarrow d_3 \dots \rightarrow d_n.$$

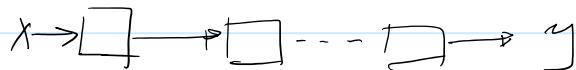
$$d_k^T A x^* = d_k^T A (\sum \alpha_i d_i)$$

$$\begin{aligned} Ax = b \\ Ax^* = b \end{aligned} \implies \Delta d_k^T b = \underline{\underline{\alpha_k}} d_k^T A d_k$$

3) Training Deep Neural Network.

① Gradient calculation. : BP

Δ ② Vanishing/Exploding gradients : BN
Batch normalization



$$y = (\pi w^k) x.$$

③ generalization NN,

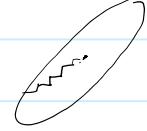
Week6

2018年2月26日 17:39

Review, Optimization methods.

$$\text{Opt. problem : } X^* = \arg \min_X f(x).$$

$$\text{GD : } X_{n+1} = X_n - \underbrace{\varepsilon \nabla f(x)}_{\nabla f(X_n)} \Big|_{x=X_n}.$$



$$\text{Newton : } X_{n+1} = X_n - \varepsilon (\nabla^2 f(x_n))^{-1} \nabla f(x_n)$$

$(\dim(X))^2$ in complexity

→ quasi newton

$$\text{Conjugate gradient} \quad d_n = -\nabla f(x_n) + \underbrace{\beta_n d_{n-1}}_{d_n}$$
$$X_{n+1} = X_n + \varepsilon d_n$$

5) Training Deep Neural Network

① Stochastic GD → SGD

② Several Neural Network problems.

Vanishing/exploding gradient : BN (Batch normalization)

SGD being too slow — momentum
adaptive learning rate

Generalization — regularization
— drop out
— L1, L2 regularization

③ Gradient Calculation, BP

① Stochastic GD:

$$h(x) = h(x|w)$$

$$w^* = \min_w f(w), \quad f(x) = E[(y - h(x|w))^2]$$

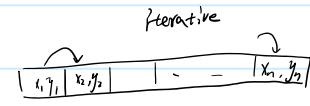
True CD:

$$w_{n+1} = w_n - \varepsilon D_n E[(y - h(x|w))^2]$$

Batch method:

$$w_{n+1} = w_n - \varepsilon D_n E[(y - h(x|w))^2]$$

$$\approx w_n - \varepsilon D_n \left(\frac{\sum_i (y_i - h(x_i|w))^2}{n} \right)$$

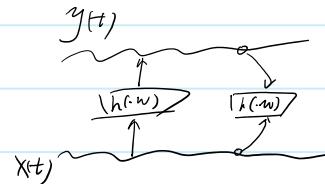


$\{(x_i, y_i)\}$: training data.

$$SGD, \quad w_{n+1} = w_n - \varepsilon D_w (y_j - h(x_j|w))^2$$

/ Randomly choose one data point at one time.

LMS/adaptive filtering

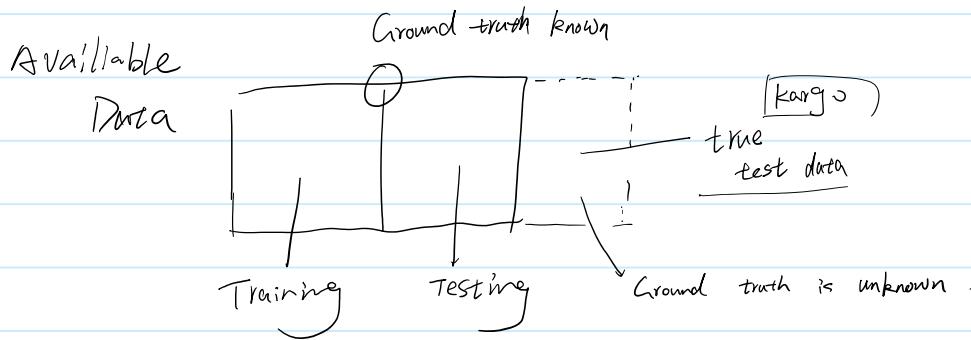


most used

mini batch SGD, random subset. / specify the size of mini batch.

$$w_{n+1} = w_n - \varepsilon D_n \sum_{k \in D_k} \frac{(y_{ik} - h(x_{ik}|w))^2}{k} \quad (n = N_n)$$

D_k = random subset of D w/ k samples.



② Several neural network problems.
 Deep)

— Proposed solutions.

<1> SGD too slow reason, GD is too slow

Proposed solution : momentum method

m_n = new direction for generating w_{n+1}

$$w_{n+1} = w_n + \varepsilon m_n.$$

$$m_n = \beta m_{n-1} - \alpha \nabla J(w) / \|w\|_2$$

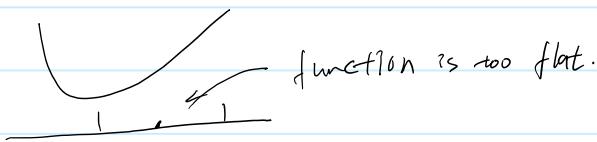
$$J(w) = E[(y - h(x/w))^2]$$

$$\approx \frac{1}{K} \sum (y_{ik} - h(x_{ik}/w))^2.$$

Conjugate gradient descent,

$$d_n = \beta_n d_{n-1} - \alpha \nabla J(w)$$

$$w_{n+1} = w_n + d_n.$$



Proposed solution. increasing learning rate

1st Step. detect flat region.

$$|\nabla f(w)| \text{ too small}$$

RMS prop

$$w_{n+1} = w_n - \frac{\alpha \nabla_w J(w_n)}{\sqrt{s_n + \epsilon}}$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} / \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ \frac{2}{4} \end{bmatrix}$$

point wise
↓

s_n becomes small when $\nabla J(w_n)$

RMS prop \rightarrow Adam

Momentum $\qquad\qquad\qquad =$

$$s_n = \beta s_{n-1} + (1-\beta) \nabla J(w_n) \cdot \nabla J(w) \Rightarrow \begin{bmatrix} 1 \\ 2 \end{bmatrix} / \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} \frac{1}{3} \\ \frac{2}{4} \end{bmatrix} \quad \beta \in (0, 1)$$

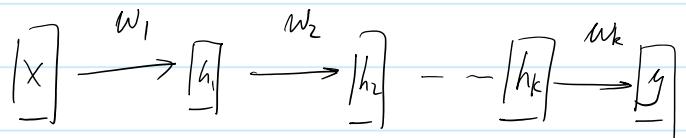
<2> Vanishing / Exploding Gradient.

$$\left(\prod_{m=1}^M w_m \right) M \text{ is large.}$$

$$\begin{aligned} w_m > 1 &\Rightarrow \text{prod} \rightarrow +\infty \\ w_m < 1 &\Rightarrow \text{prod} \rightarrow 0 \end{aligned}$$

finite product, most $w_m = 1$

Neural networks



$$h_1 = \phi(w_1, x)$$

Normalize $w_k \cdot h_k$ around 1

$$h_2 = \phi(w_2, h_1) = \phi(w_2 : \phi(w_1, x))$$

1

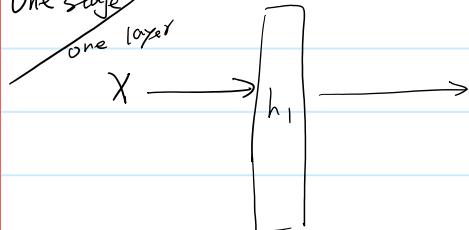
$$y = \phi(w_k, h_k) = \phi(-\phi(w_k, h_k))$$

$$\phi = I \quad \Rightarrow \quad y = (\pi w_k)_x$$

\searrow linear case.

Batch
normalization

One stage



$$h_1 = \phi(w, x)$$

$$= \emptyset(u) \quad u = w_1 \cdot x.$$

Want u_i to have a $N(0, 1)$

Gaussian $\{E^{(n)}_{\text{Var}}\}$

$$t = \frac{u_i - E[u]}{\text{Var}[u]} \leftarrow \begin{array}{l} \text{constant} \\ \text{calculate using} \\ \text{mini-batch} \end{array}$$

$$t \sim N(0, 1)$$

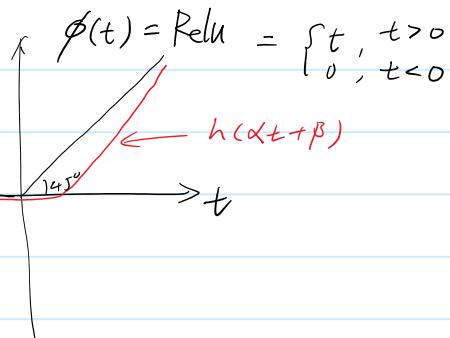
Replace $b_1 \phi(w, x) = \phi(w)$

$$\text{by } h_1 = \phi(\alpha t + \beta)$$

undated parameter
to be trained.

$$t \sim N(0, 1)$$

$$\alpha t + \beta \sim N(\beta, \alpha^2)$$



w/o normalization $\alpha, \beta \rightarrow w_k$.

w normalization need α, β

At Step n

$$\alpha_n, \beta_n \xrightarrow{w_n} t_n \text{ (through normalization)}$$

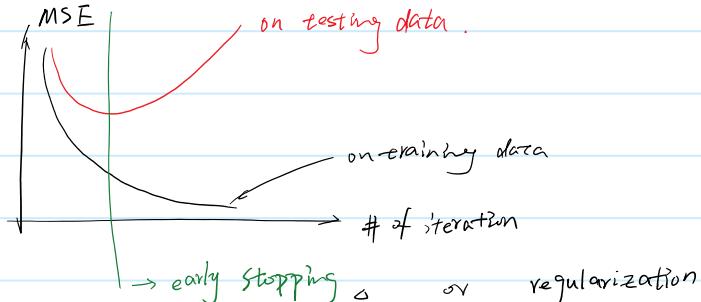
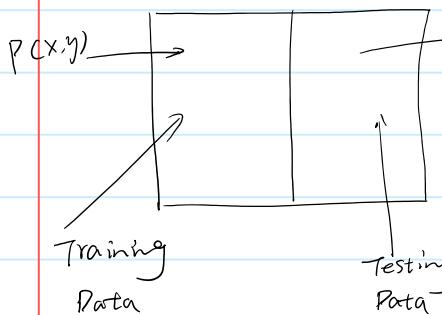
$$\Rightarrow h_i = \phi(\alpha_n t_n + \beta_n)$$

→ calculate Gradient wrt w_n, α_n, β_n .

$$\rightarrow \alpha_{n+1}, \beta_{n+1}, w_{n+1}$$

→ Step $n+1$

(5) Generalization.



$$MSE = E[(y - h(x|w))^2]$$

On training data.

$$= E_p[(y - h(x|w))^2] \approx \frac{1}{n} \sum (y_i - h(x_i|w))^2$$

On testing data

$$= E_p[(y - h(x|w))^2] \approx \frac{1}{n} \sum (y_j - h(x_j|w))^2$$

Proposed solution ① early stopping
② regularization

example, 10,000 samples.

10^6 weights

under constrained problem.

→ introduce constraint, most weights equals to zero

10 non zero weights

→ $\binom{10}{10,000}$ ⇒ too many comb.
to check.

Solution, regularization.

b₁ regularization.

$$\begin{bmatrix} 1 \\ 2 \\ \vdots \\ 10 \end{bmatrix} \Rightarrow \text{non sparse}$$

$$\nabla J(\omega) = E[(y - h(x|\omega))^2]$$

$$\rightarrow \nabla J(\omega) + \lambda \|\omega\|_1$$

$$\begin{bmatrix} 1 \\ 2 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \text{sparse} \quad \text{most are 0.}$$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad \|x\|_1 = \sum_{i=1}^m |x_i|.$$

⇒ sparse solution

$$b_2 \rightarrow \nabla J(\omega) + \lambda \|\omega\|_2$$

$$\|x\|_2 = \sqrt{\sum x_i^2} \quad x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

⇒ small weights

Week7

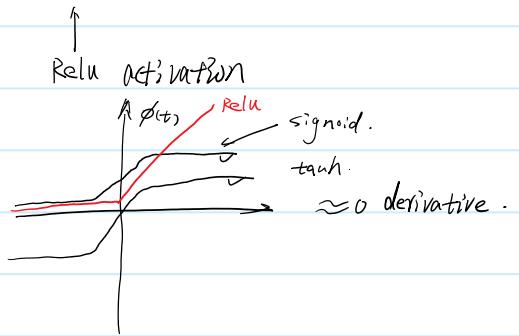
2018年3月5日 19:14

Training Deep Neural Networks

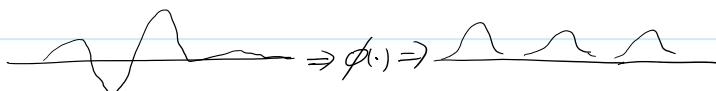
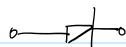
— calculate gradient, BP

— exploding/vanishing gradient, Batch normalization.

Selecting initial value — He initialization



Rectification



$\phi(\alpha t + \beta)$

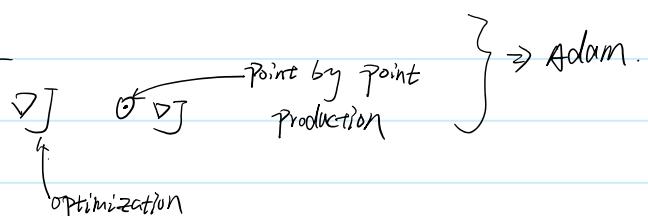
normalized
 $N(0, 1)$

— GD too slow

— momentum

similar to CGM

— adaptive learning rate,

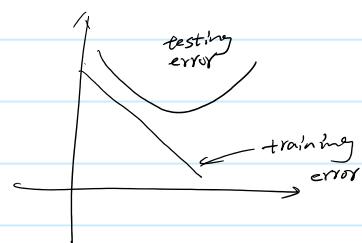


— Regularization,

training error, already small & decrease

testing error,

& increase



l_2, l_1 regularization

early stop

drop off / randomly drop off a set of nodes.

$$\text{rate} = \lambda$$

train the rest

Gradient calculation

Review the chain rule.

$$f(x) = y \quad \frac{dy}{dt} = \frac{dg}{dx} \cdot \frac{dx}{dt} = f' \cdot g'$$

$$x = g(t)$$

$$t \xrightarrow{g} x \xrightarrow{f} y$$

$$(t, s) \xrightarrow{g} (x, y) \xrightarrow{f} (u, v)$$

$$\frac{\partial u}{\partial t} = \frac{\partial u}{\partial x} \cdot \frac{\partial x}{\partial t} + \frac{\partial u}{\partial y} \cdot \frac{\partial y}{\partial t}$$

$$(x_1, \dots, x_N) \rightarrow (v_1, \dots, v_I) \rightarrow (u_1, \dots, u_J) \rightarrow (y_1, \dots, y_K)$$

input

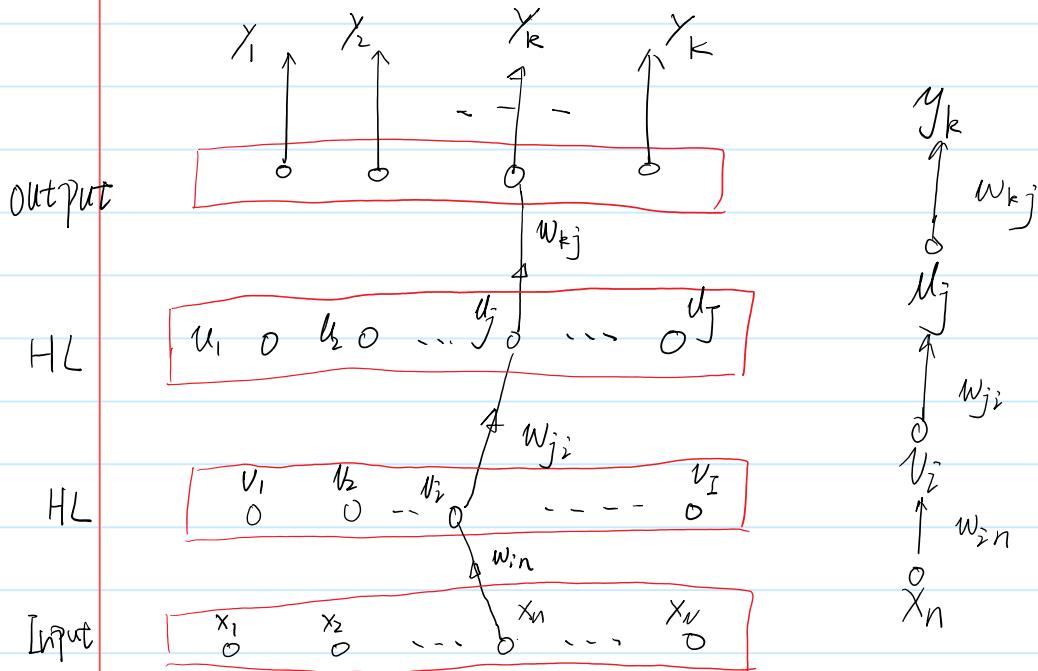
HD 1

HD 2

output

$$\frac{\partial y_k}{\partial v_i} = \sum_j \frac{\partial y_k}{\partial u_j} \cdot \frac{\partial u_j}{\partial v_i}$$

$$\frac{\partial y_k}{\partial x_n} = \sum_{i,j} \frac{\partial y_k}{\partial v_j} \cdot \frac{\partial u_j}{\partial v_i} \cdot \frac{\partial u_i}{\partial x_n}$$



Training

$$J(w) = \sum \frac{1}{2} \|y - t\|^2$$

(t, x)
target output
input

$$J(w) = \sum \frac{1}{2} \|NN(x_i, w) - y_i\|^2.$$

$D = \{(x_i, y_i)\}$

Gradient calculation:

$$\frac{\partial J}{\partial w} = \left\{ \text{all partials of } J \text{ wrt weight} \right\}$$

$\overbrace{}$
tensor

$$J(W) = \frac{1}{2} \|y(x) - t\|^2. \quad (\text{one sample mini batch.})$$

$$= \frac{1}{2} \sum_k (y_k(x) - t_k)^2.$$

↓
 input sample.
 Variable.

↓
 target value.
 $y_k(x) = y_k(x, w).$

$$= \frac{1}{2} \sum_k (y_k(x, w) - t_k)^2.$$

$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial}{\partial w_{kj}} \frac{1}{2} \sum_k (y_k - t_k)^2$$

↓
 variable
 ↓
 target value.

$$= \underbrace{\sum_k (y_k - t_k)}_{\text{error } \delta_k} \cdot \frac{\partial}{\partial w_{kj}} (y_k - t_k)$$

$$= \sum_k \delta_k \frac{\partial y_k}{\partial w_{kj}} = \delta_k \frac{\partial y_k}{\partial w_{kj}} = \underbrace{\delta_k}_{\text{error}} \underbrace{\phi'(\sum_j w_{kj} u_j)}_{\text{activation function}} \cdot u_j$$

$$y_k = \phi(\sum_j w_{kj} u_j)$$

$$\frac{\partial y_k}{\partial w_{kj}} = \phi'(\sum_j w_{kj} u_j) \cdot \frac{\partial}{\partial w_{kj}} (\sum_j w_{kj} u_j)$$

$$= \phi'(\quad) u_j$$

When ϕ is linear $\phi'() = 1$

$$\boxed{\frac{\partial J}{\partial w_{kj}} = \delta_k u_j}$$

$$\frac{\partial J}{\partial w_{ji}} = \frac{2}{\partial w_{ji}} \times \frac{1}{2} \sum_k (y_k - t_k)^2$$

$$= \sum_k (y_k - t_k) \frac{\partial y_k}{\partial w_{ji}}$$

$$= \sum_k \delta_k \frac{\partial y_k}{\partial w_{ji}}$$

$$\frac{\partial y_k}{\partial w_{ji}} = \frac{\partial y_k}{\partial u_j} \cdot \frac{\partial u_j}{\partial w_{ji}}$$

$$y_k = \phi(\sum_j w_{kj} u_j)$$

$$\frac{\partial y_k}{\partial u_j} = \phi'(\sum_j w_{kj} u_j) \cdot w_{kj}$$

$$= \phi'(\quad) \cdot w_{kj}$$

$$\frac{\partial J}{\partial w_{ji}} = \frac{\partial}{\partial w_{ji}} \phi\left(\sum_i w_{ji} v_i\right)$$

$$= \underline{\phi'\left(\sum_i w_{ji} v_i\right) \cdot v_i}$$

$$\begin{aligned} \frac{\partial J}{\partial w_{ji}} &= \sum_k \delta_k \frac{\partial y_k}{\partial w_{ji}} \\ &= \sum_k \delta_k \phi'\left(\sum_j w_{kj} u_j\right) \cdot w_{kj} \cdot \phi'\left(\sum_i w_{ji} v_i\right) v_i \\ &= \underbrace{\sum_k \delta_k \phi'\left(\sum_j w_{kj} u_j\right) \cdot w_{kj}}_{\text{linear combination error at layer above.}} \cdot \phi'\left(\sum_i w_{ji} v_i\right) v_i \\ &\xrightarrow{\substack{\text{error} \\ \text{back-propagation}}} \delta_j \end{aligned}$$

$$= \delta_j \cdot \phi'\left(\sum_i w_{ji} v_i\right) v_i \Rightarrow \boxed{\frac{\partial J}{\partial w_{ji}} = \delta_j \phi' v_i}$$

$$\frac{\partial J}{\partial w_{kj}} = \delta_k y'_k u_j$$

$$\delta_j = \sum_k \delta_k y'_k w_{kj}$$

$$\frac{\partial J}{\partial w_{ji}} = \delta_j u'_j v_i$$

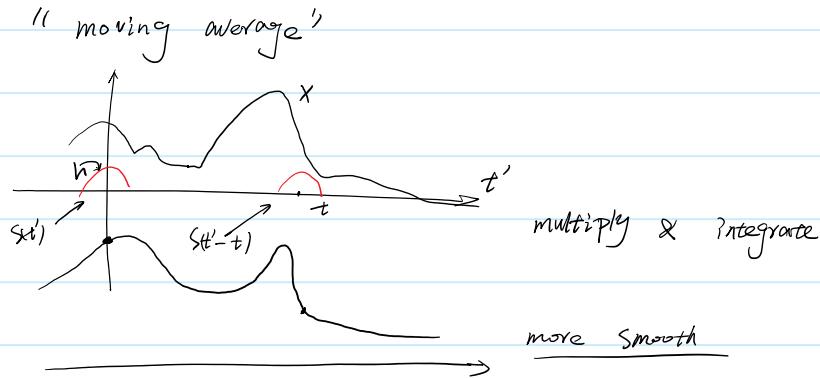
$$\delta_j = \sum_k y'_k w_{kj} \delta_k$$

$$\frac{\partial J}{\partial w_{in}} = \delta_i v'_i x_n$$

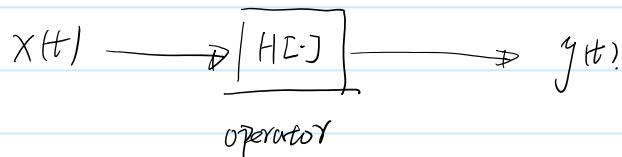
$$\delta_i = \sum_j \delta_j u_j w_{ji}$$

4. Convolutional Neural Networks (CNNs)

1) Convolution.



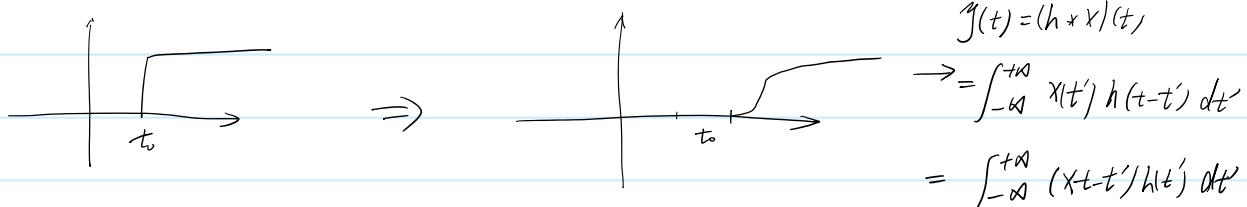
① linear time invariant systems (LTI).
(shift)



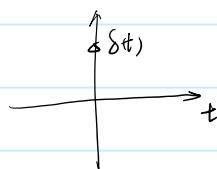
linear: $H[a x_1(t) + b x_2(t)] = a H[x_1(t)] + b H[x_2(t)]$



② convolution (LTI system)



$H[\cdot]$, LSI complete det. by $H[\delta(t)]$



$\delta(t)=0, t \neq 0$

$\int \delta(t) dt = 1$

$H[\delta(t)] = h(t)$, impulse response

$$y(t) = \int_{-\infty}^{+\infty} x(t') h(t - t') dt'$$

flip and move

$$= \int_{-\infty}^{+\infty} x(t-t') h(t') dt'.$$

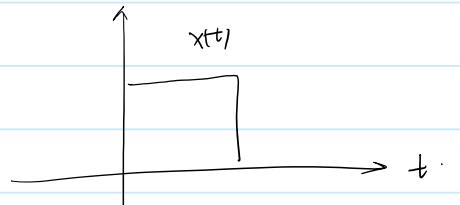
$$x(t) \rightarrow X(\omega)$$

$$h(t) \rightarrow H(\omega)$$

$$Y(\omega) = H(\omega) X(\omega).$$

$$y(t) \rightarrow Y(\omega)$$

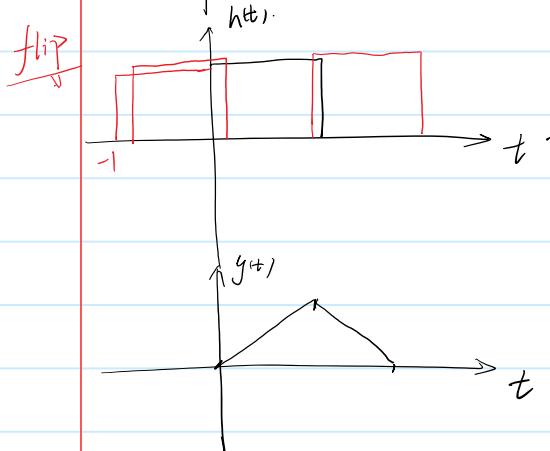
filter, h(t), H(ω)



$$\text{Support}(x(t)) = 1$$

$$(h(t)) = 1$$

{ the length of interval
on which it is not zero



$$y(t) = (x * h)(t)$$

$$\text{Support}(y(t)) = |+| = 2.$$

$$\text{Support}(h_1 * h_2 * h_3 * h_4) = \sum_{i=1}^4 \text{Support}_i$$

Week8

2018年3月12日 15:38

Review

Convolution \leftrightarrow moving average.

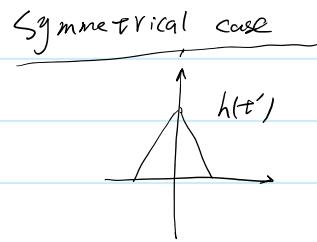
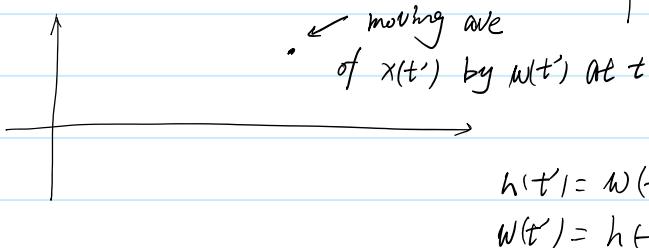
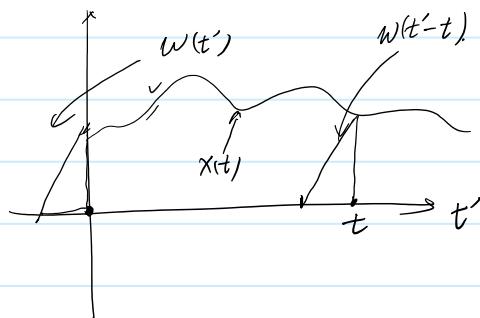
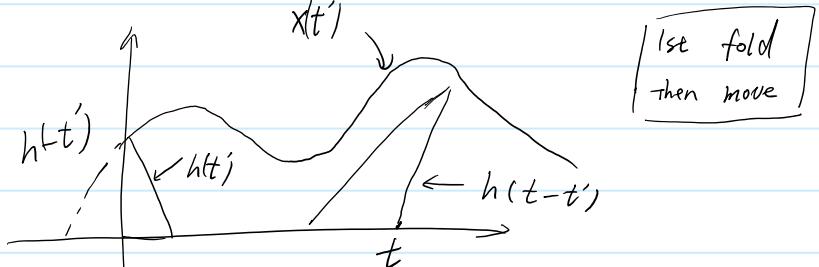
$$x(t) \downarrow \quad h(t) \rightarrow \oplus \Rightarrow y(t) \quad (x * h)(t) = \int_{-\infty}^{+\infty} x(t') h(t-t') dt'$$

or

$$\int_{-\infty}^{+\infty} x(t-t') h(t') dt'$$

moving average.

$$y(x(t) w/w(t)) = \int_{-\infty}^{+\infty} x(t) w(t-t') dt' \\ = y(t)$$



④ Discret Case $\not\rightarrow$ in computer.

$$x(t) \rightarrow X[n] \quad (x_n)$$

$$h(t) \rightarrow h[n]$$

$$y(t) \rightarrow y[n]$$

Discrete convolution

$$y[n] = \sum_m x(m) h(n-m)$$

moving ave

$$y[n] = \sum_m x(m) w(m-n)$$

$x[0], x[1], \dots, x[N-1]$ input sequence.

$h[0], h[1], \dots, h[p-1] \quad p \ll N$

$$\begin{bmatrix} y[0] \\ y[1] \\ \vdots \\ y[N-1] \end{bmatrix} = \begin{bmatrix} h[0] & 0 & \cdots & 0 \\ h[1] & h[0] & 0 & \cdots & 0 \\ h[2] & h[1] & h[0] & \ddots & 0 \\ \vdots & & & & \\ h[p-1] & \cdots & h[0] & 0 & \cdots & 0 \\ 0 & h[p-1] & \cdots & h[0] & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[N-1] \end{bmatrix}$$



△ matrix very sparse.

△ row circular shifted.

⑤ Subsample conv result.

Subsample
by 2.

$$\begin{array}{c} y[0] \checkmark \\ y[1] \\ \cancel{y[2]} \\ \cancel{y[3]} \\ y[4] \checkmark \\ \vdots \\ y[N-1] \end{array} = \begin{bmatrix} h[0] & 0 & \cdots & 0 \\ h[1] & h[0] & 0 & \cdots & 0 \\ h[2] & h[1] & h[0] & \ddots & 0 \\ \vdots & & & & \\ h[p-1] & \cdots & h[0] & 0 & \cdots & 0 \\ 0 & h[p-1] & \cdots & h[0] & 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[N-1] \end{bmatrix}$$

↓
Circular shift by 2

⇒ stride

⑥ Extension to 2D

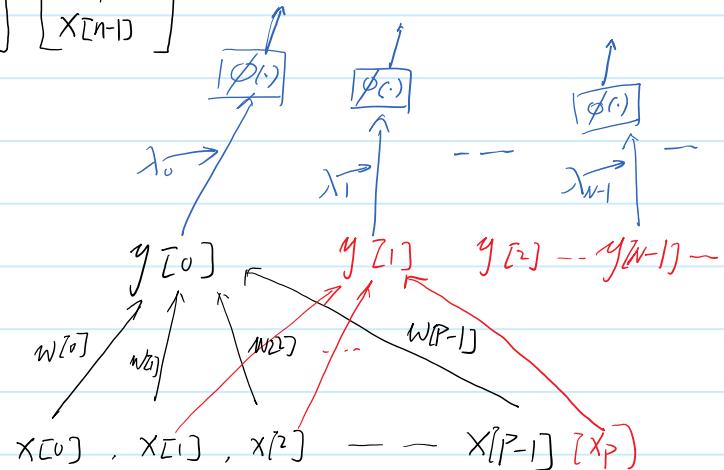
$$y[m, n] = \sum_{m', n'} x[m', n'] h[m-m', n-n']$$

⑦ Neural network implementation of conv./moving ave.

$$\begin{bmatrix} y[0] \\ y[1] \\ y[2] \\ y[3] \\ y[4] \\ \vdots \end{bmatrix} = \begin{bmatrix} w[0] & w[1] & \cdots & w[P-1] & 0 & \cdots & 0 \\ 0 & w[0] & \cdots & \cdots & w[P-1] & 0 & \cdots & 0 \\ 0 & 0 & w[0] & \cdots & \cdots & w[P-1] & 0 & \cdots & 0 \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[n-1] \end{bmatrix}$$

$$w[m] = w[0], w[1], \dots, w[P-1]$$

$$P \ll N$$



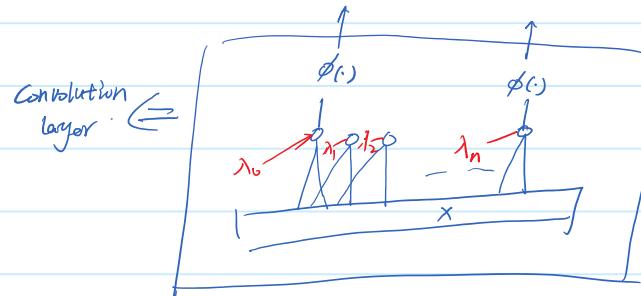
Weights for entire layer.

$$h = \phi(w * x + \lambda)$$

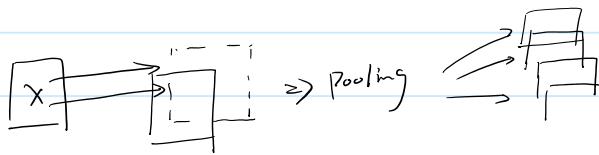
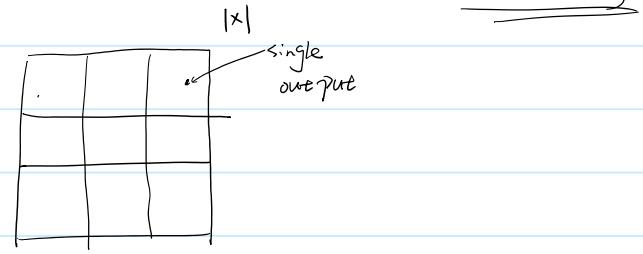
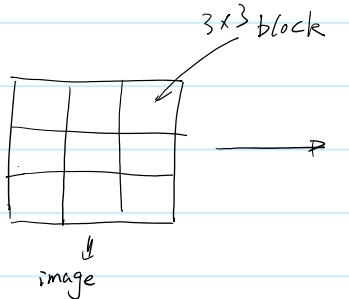
Vector applied component wise (threshold & nonlinearity)

↓

to have nonlinearity complexity in function



⑧ Pooling



feature filtering
 $2 \times 4 = 8$ types

(different conv kernel)
ex. different orientations

⑨ Notable CNNs

LeNet 5

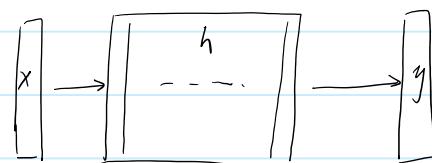
AlexNet

Google Net

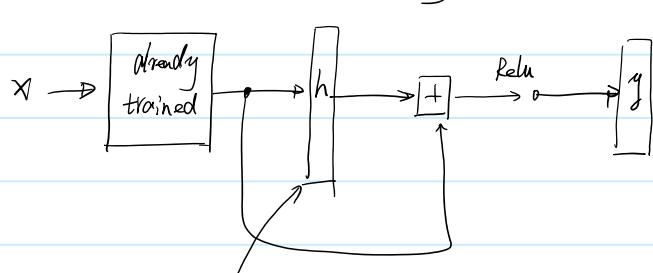
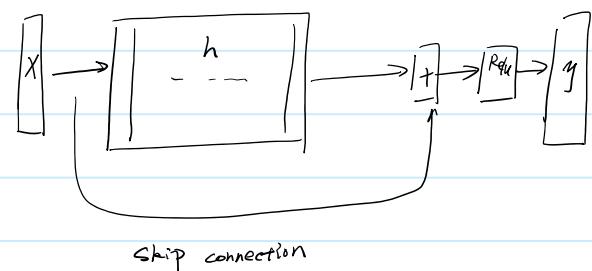
ResNet 152 layers

ResNet

conventional network



ResNet



initialized $\sim N(0, \sigma^2)$, $\sigma^2 \ll 1$.

training starting from already trained; only accept things that makes it better

Activation function

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU) [2]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

5

2018年3月26日 16:58

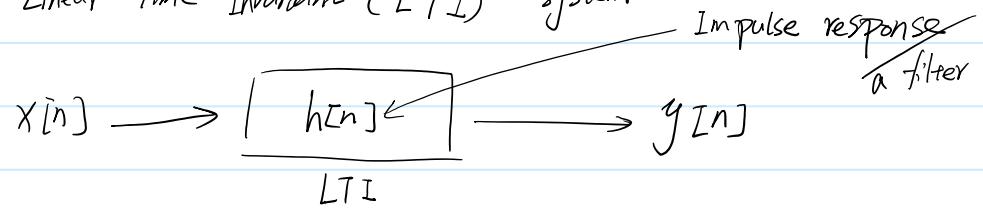
Week9

2018年3月26日 17:38

5. Recurrent Neural Network (RNNs)

1) FIR and IIR filters

① Recall Linear Time Invariant (LTI) system



$$y[n] = (h * x)[n]$$

② FIR system (Finite impulse response)

$h[n]$ has finite support

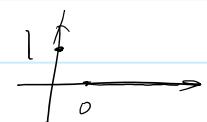
$$y[n] = (h * x)[n] = \text{local average.}$$

example: $y[n] = \frac{1}{3} (x[n] + x[n-1] + x[n-2])$

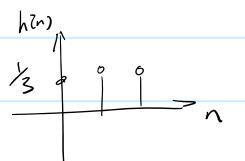
$$h[n] = ?$$

$$h[n] = \text{impulse response} = \text{output} (\text{input} = \delta(n))$$

$$= \frac{1}{3} [x[n] + x[n-1] + x[n-2]] / \quad \delta(n) = \begin{cases} 0, & n \neq 0 \\ 1, & n = 0 \end{cases}$$



$$= \frac{1}{3} [\delta(n) + \delta(n-1) + \delta(n-2)]$$



⑧ IIR (infinite impulse response)

example: $y[n] = \alpha y[n-1] + (1-\alpha)x[n]$

\downarrow
memory

$h[n] = ?$ Z transform

recursive equation

$$\underline{y[-1] = 0}$$

$h[n] = \text{output} (\text{input } x[n] = \delta[n])$

$$h[0] = y[0] = \alpha y[-1] + (1-\alpha)x[0] \quad | \quad x[0] = \delta[0]$$

\uparrow
0 \uparrow
1

$$= 1 - \alpha$$

$$h[1] = y[1] = \alpha y[0] + (1-\alpha)x[1] \quad | \quad x[1] = \delta[1]$$

\uparrow
 α \uparrow
0

$$= \alpha(1-\alpha)$$

$$h[2] = y[2] = \alpha y[1] + (1-\alpha)x[2] \quad | \quad x[2] = \delta[2]$$

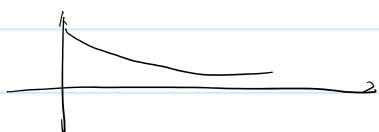
\uparrow
 $\alpha(1-\alpha)$ \uparrow
0

$$= \alpha^2(1-\alpha)$$

$$h[n] = (1-\alpha) \alpha^n, \quad n = 0, 1, 2, \dots$$

Support [$h[n]$] = $+\infty$

$$y[n] = [h * x][n]$$



$| \alpha | < 1$ for system to be stable

2) Relation to Neural Network.

FIR filter \longrightarrow CNN

IIR filter \longrightarrow RNN

$$\text{CNN layer} \quad x \rightarrow \underset{\substack{\text{CNN} \\ \text{layer}}}{\text{CNN}} \rightarrow y$$

$$y = \phi(\underbrace{(x * h)}_{\substack{\text{Small kernel}}} + \lambda)$$

$$= \phi(\underbrace{(W_x x)}_{\substack{\equiv \\ \equiv}} + \lambda)$$

$$\boxed{=} \quad \boxed{=}$$

① RNN, Simple RNN

$$y[n] = \phi(W_x x[n] + W_y y[n-1] + \lambda)$$

Vector

non-linear
applied
component wise

$$y[n] = \phi(\alpha y[n-1] + (1-\alpha) x[n] + \lambda)$$

Vector

Scalar

Scalar

Scalar

graph

advantage,

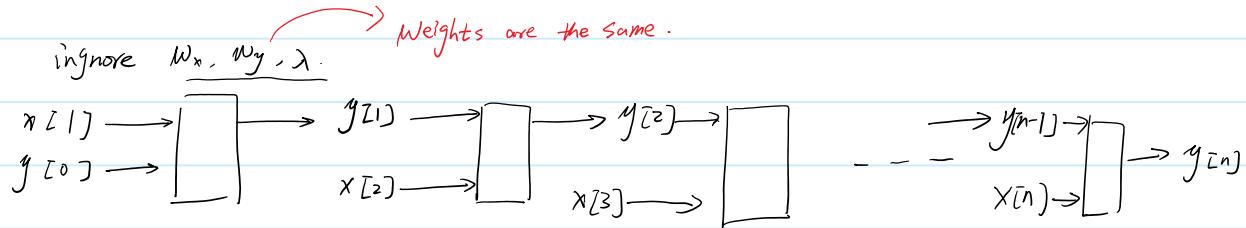
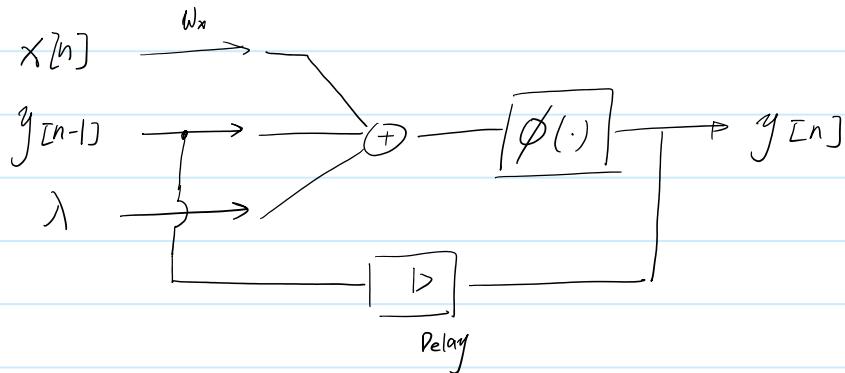
① non-linear

② system with long memory (infinite window/kernel size)

③ dimension of input not need to be constant.

obvious direction in time

(2) Roll-out representation



Input sequence of length n .

Output sequence of length n .

(3) Training in RNN

Typically training Data.

$$D = \left\{ \left(\overset{\text{Batch}}{\underset{1}{\overset{n}{\cdots}}} \left(x_{[1]}^1, y_{[0]}^1 \right), \left(x_{[2]}^1, y_{[1]}^1 \right) \dots \left(x_{[n]}^1, y_{[n-1]}^1 \right), y_{[n]}^1 \right) \right\}$$

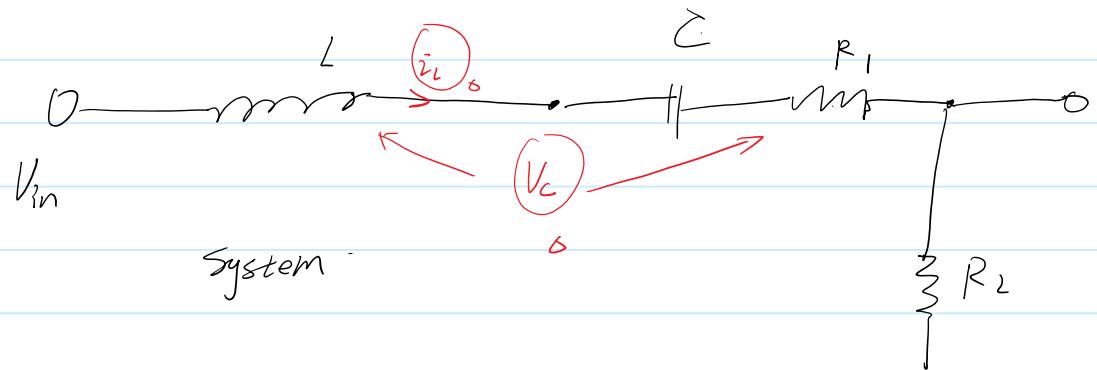
$$\left(x_{[1]}^2, y_{[0]}^2 \right), \left(x_{[2]}^2, y_{[1]}^2 \right) \dots \left(x_{[n]}^2, y_{[n-1]}^2 \right), y_{[n]}^2$$

↑
|
|
|

4) More general RNN.

Simple RNN, extension of $y[n] = \alpha y[n-1] + (1-\alpha)x[n]$

more general RNN, (1) State space representation.



Linear case

State variables $x[\lambda]$

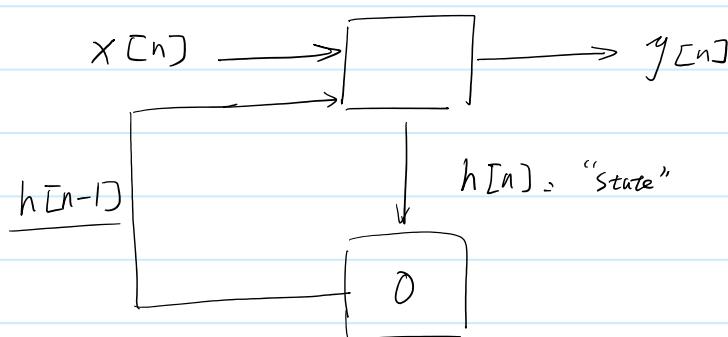
state evolution

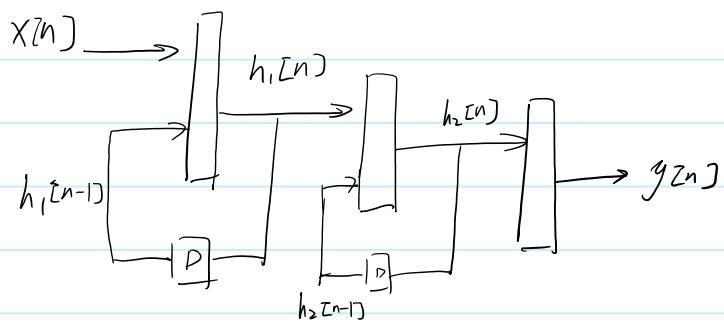
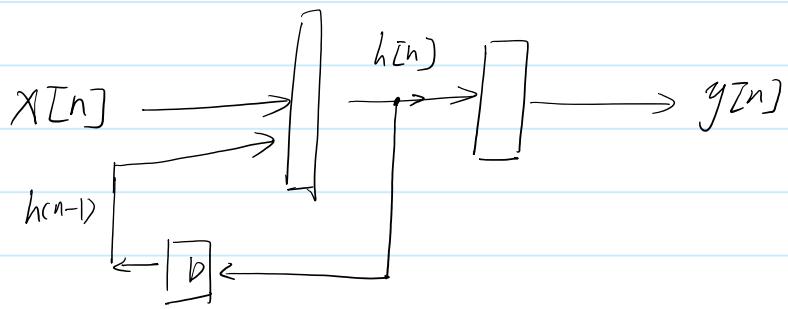
$$x[n+1] = Ax[n] + w[n]$$

$$y[n] = Hx[n] + v[n]$$

observational ↑
 may or may not be
 observational

$$\Rightarrow E[x[n] | y[n], y[n-1], \dots, y[0]]$$





Week10

2018年4月2日 17:45

clustering (Grouping)

1) What?

① Example, cluster/group gradient.

W. HF \rightarrow EF.

HF: heart failure.

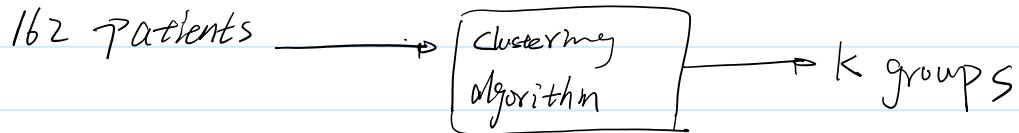
EF: ejection fraction

< 50% abnormal

P: preserved.

normal

65 features / patient $\longrightarrow x_i$



② Data Grouping

w/in a group: similar

btw group: dis-similar

\Rightarrow unsupervised classification

③ Key to clustering

- distance measure.

- data model

2) distance measures

I.

real Valued data

$$x \in \mathbb{R}^n, y \in \mathbb{R}^n$$

$$d(x, y) = \|x - y\|. \quad (\text{Euclidean})$$

normalized $d(x, y) = \|x - y\|_{\Sigma} \leftarrow \text{matrix}$

$$= \sqrt{(x-y)^T \Sigma^{-1} (x-y)}$$

II.

$$d_E(x, y) < d_E(x, z)$$

$$d(x, y) > d(x, z)$$

III. Categorical data

Example.

$$x = (\text{yes}, \text{yes}) \rightarrow (0, 0)$$

$$y = (\text{no}, \text{no}) \rightarrow (1, 1)$$

$$d(x, y) = d_H(x, y) = 2$$

Hamming distance.

categorical but not binary.

Red $\rightarrow 001$ Blue $\rightarrow 010$ Green $\rightarrow 100$

$$x = (\text{yes}, \text{red}) \rightarrow (0, 001)$$

$$y = (\text{no}, \text{green}) \rightarrow (1, 100)$$

$$z = (\text{no}, \text{blue}) \rightarrow (1, 010)$$

$$\underline{d_H(x, y) = 3.}$$

Problem

$$X = (x_1, x_2 \dots x_n)$$

x_i = Categorical Variable.

w/ 10^7 categories

One-hot rep of x_i , 10^7 long

When does this happen? x_i is a word.

is a comb. of cate variables.

Word-Embedding entity 10^7 dim. one-hot vector \rightarrow m dim real vector

- $m \ll 10^7$
- if x and y are similar

$$\|x_i^e - y_i^e\| \text{ small.}$$

Rough Idea: Find the embedding in a classification task.

by a neural network

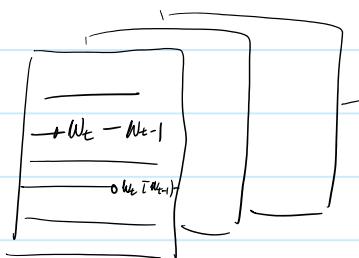
$$\underline{N\text{-Gram}} : P(w_t | w_{t-1} \dots w_{t-(n-1)}) = P(w_t | w^{t-1})$$

assume to be true prob. distribution for words.

$$g(w_t | w_{t-1}) = \text{estimate of } P(w_t | w^{t-1})$$

= output a neural network

Training data: document / articles.



dog
cat
:
 w_t
 w_{t-1}

$$P(\text{dog} | w_{t-1}) = 1.$$

$$P(\text{cat} | w_{t-1}) = 0$$

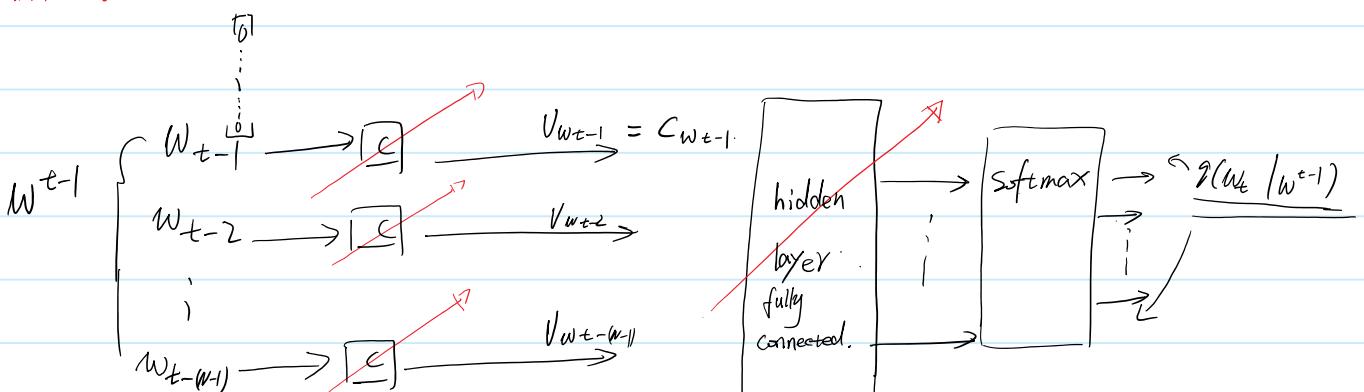
optimization function. cross entropy

$$\max \sum_{w \in t} - \sum_{w' \in t} P(w_t' | w^{t-1}) \log g(w_t' | w^{t-1})$$

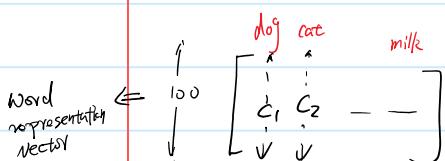
all possible words either one or zero.

all training data point.

All point in mini-batch.



real vector for each word.



things change

$$[c_1 \dots 0 \dots] = c_n$$

3) Hierarchical clustering

bottom up.

Distance-based.

no model.

(Medical)

Adv.

Simple.

works well when there is
hierarchy in data

problems,

where to cut

4) K-means LBG, VQ

K , the number of clusters. Given

initialization, start w/ K cluster centers. $C_k^{(0)}$, $k=1, 2, \dots, K$.

classification, assign x_i to $C_k^{(n)}$ \leftarrow n^{th} iteration $n=0, 1, 2, \dots$

$$\text{if } d(x_i, C_k^{(n)}) = \min_{k'} d(x_i, C_{k'}^{(n)})$$

min distance classification.

update:

C_k^{n+1} = average of all points in the k^{th} cluster.

check:

$$\|C^{n+1} - C^n\| < \varepsilon$$

$$C = (C_1, C_2, \dots, C_K)^T$$

yes \rightarrow stop

No \rightarrow go to classification

Adv: simple.
 works well if
 a good K
 the data is spherical.

Problems,

K is not known.

If data is not spherical.

initialization, needs good initialization.

solution

→ modeling

Gaussian mixture

1) Info theoretic criteria
 (AIC, BIC, XIC) .

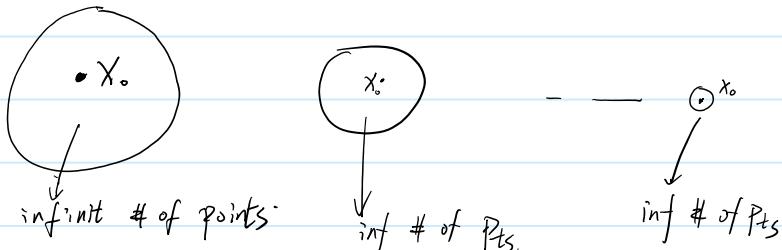
2) non parametric Bayesian.

3) EM Bayesian.

5) Density - Distance clustering

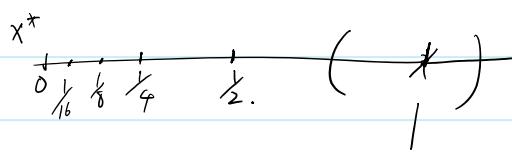
Basic idea: math analysis.

A cluster point.



$$D = \{x_i \mid x_i = 2^{-i}\} \quad i \geq 0, \text{ integer}$$

cluster point for D : $x^* = 0$



Week11

Monday, April 9, 2018 5:44 PM

Review, clustering Algorithm.

- hierarchical
 - k-means
 - density-distance
- = every data point
→ a cluster.

fuzzy - k-means

Model-based clustering

- A pdf for data.

- MLE / MAP → $\hat{\mu}_s$

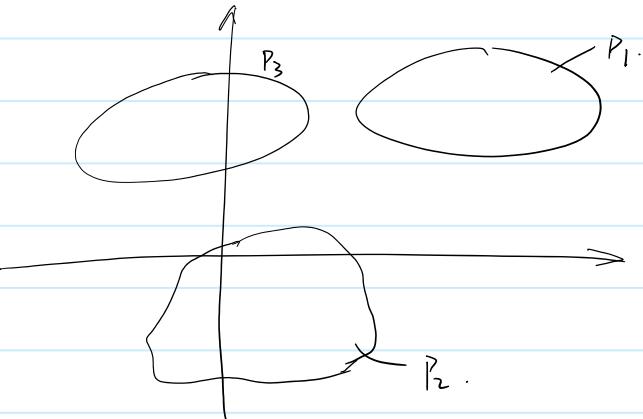
EM algorithm

6) Model based clustering & EM Algorithm.

① mixture model

$$P(x) = \sum_{k=1}^K \pi_k P_k(x)$$

$$\pi_k > 0, \sum \pi_k = 1$$

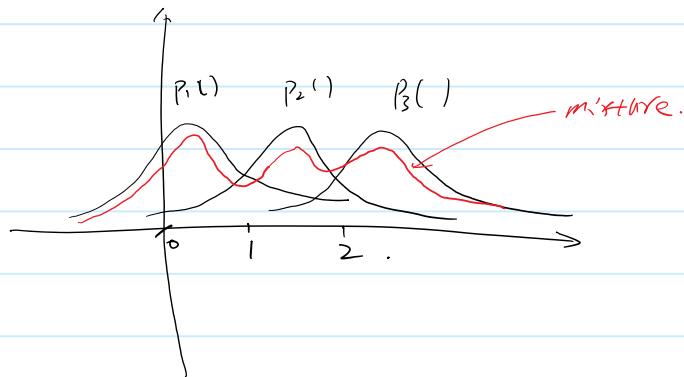


Example, Gaussian mixture

$$P(x) = \sum_{k=1}^K \pi_k P_k(x | m_k, \sigma_k^2)$$

$$\text{e.g. } K=3, \pi_1 = 0.2, \pi_2 = 0.4, \pi_3 = 0.4$$

$$\sigma_1 = \sigma_2 = \sigma_3 = 1.$$



Relation to clustering

Δ fix k.

Data $D = \{y_1, y_2, \dots, y_n\}$ i.i.d.

$$\text{model: } p(y) = \sum_{k=1}^K \pi_k P_k(y)$$

$$= \sum_{k=1}^K \pi_k P_k(y|\theta_k) = p(y|\theta)$$

$\theta = (\{\pi_k\}, \{\theta_k\})$.

Training $\triangleright \hat{\theta} \rightarrow P[y_i \in \text{clustering} | y_i, \hat{\theta}]$

Δ finding k

- info. theoretic criteria.
- NPB (Non-Parametric Bayesian).
- ✓ - as a part MLE/MAP

② MLE for clustering

clustering $\triangleright \hat{\theta}$

$$= \underset{\theta}{\operatorname{argmax}} \log P(D|\theta)$$

$$\text{MLE: } \theta = \underset{\theta}{\operatorname{argmax}} \log P(y_1, y_2, \dots, y_n|\theta) = \underset{\theta}{\operatorname{argmax}} \log P(D|\theta)$$

$$\log P(D|\theta) = \log P(y_1, y_2, \dots, y_n | \theta)$$

iid.

$$= \sum_{i=1}^n \log P(y_i | \theta)$$

$$= \sum \log \left(\sum_{k=1}^K \pi_k P_k(y_i | \theta_k) \right)$$

m_k, Σ_k^2
 $\theta = (\{\pi_k\}, \{m_k, \Sigma_k\})$

$$= \sum_{i=1}^n \log \left(\sum_{k=1}^K \pi_k (2\pi)^{-\frac{n}{2}} |\Sigma_k|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} \|y_i - m_k\|^2 \Sigma_k^{-1} \right] \right)$$

highly non-linear / non-quadratic

constraint: $\pi > 0, \sum \pi_k = 1$

→ difficult to solve

⑧ The EM (expectation – maximization) Algorithm.

Basic idea:

Auxiliary variable

$$P(y) = \int P(y|z) P(z) dz.$$

$$\text{Data} = \{y_1, y_2, y_3, \dots, y_n\} = y$$

↑ complex
↑
not complex
as

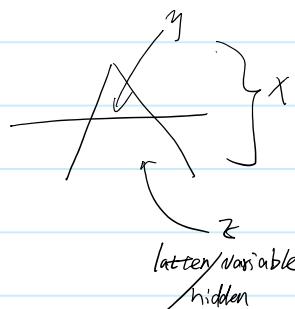
Introduce $x \rightarrow$ complete data.

→ Topics Model

y_z incomplete data.

- y incomplete data

$$y = H(x)$$



- y_z complete data. $= x$.

$$H \begin{pmatrix} y \\ z \end{pmatrix} = y$$

MLE, $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log(y|\theta)$ Difficult if done directly.

EM, iterative b/w E and M steps.

Start, initial guess — $\theta^{(0)}$

E-step, generate function

$$\tilde{\Phi}(\theta | \theta_f^{(P)}) = E \left[\log P(x|\theta) | y, \theta^{(P)} \right] \text{ current estimate}$$

$$(\approx \log P(y|\theta))$$

Good current appr. to $\log P(y|\theta)$.

M-step,

$$\theta^{(P+1)} = \underset{\theta}{\operatorname{argmax}} \tilde{\Phi}(\theta | \theta^{(P)})$$

convergence: $\log P(y|\theta^{(P+1)}) \geq \log P(y|\theta^{(P)})$

convergence to local max.

Can be highly dependent on initial guess.

(*) Applied to clustering

- Find/Def complete data.
- E-step can be difficult.
- M-step easy

Incomplete data. $\mathbf{y} = \{y_1, y_2, \dots, y_n\}_{(i, i, d)}$

complete data. $\mathbf{x} = \{y, z\}$

$$P(y) = \sum_k \pi_k P_k(y)$$

Gaussian

$$\mathbf{z} = \{z_1, z_2, \dots, z_N\}$$

= indicator vectors

$z_i = k$, if $y \in$ cluster k .

z_i , binary indicator vector

Identity

$$\downarrow$$

$$\mathbf{y} = (\mathbf{I}, \mathbf{0}) \begin{pmatrix} \mathbf{y} \\ \mathbf{z} \end{pmatrix}$$

$$= [0, 0, 0, 1, \dots, 0] \quad \text{If } y_i \in \text{cluster } k.$$

\downarrow

= "one-hot" vector

E-step,

$$\hat{\Phi}(\theta | \theta^{(T)})$$

$$\log p(x|\theta) = \log P(y, z|\theta)$$

$$= \log P(y|z) + \log P(z)$$

$$= \sum \log P(y_i|z_i) + \sum \log P(z_i).$$

$$\log P(z_i)$$

$$= - \sum \log P(y_i|z_i) + V_{\pi}^T \sum z_i$$

$$P(z_i) = \begin{cases} \pi_1, z_i = 1 = e_1 = [1, 0, 0, \dots, 0] \\ \pi_2, z_i = 2 = e_2 = [0, 1, 0, \dots, 0] \\ \vdots \\ \pi_k, z_i = k = e_k = [0, 0, 0, \dots, 1] \end{cases}$$

$$\vec{z} = [\log \pi_1, \log \pi_2, \dots, \log \pi_k]^T z_i$$

$$= V(\pi)^T z_i$$

$$\pi = (\pi_1, \pi_2, \dots, \pi_k)$$

$$V_{\pi} = \begin{bmatrix} \log \pi_1 \\ \vdots \\ \log \pi_k \end{bmatrix}$$

$$\log P(y_i | z_i) = [\log P(y_i | e_1), \dots, \log P(y_i | e_k)]_{z_i}.$$

$$= U^T(y_i) z_i.$$

$$U^T(y_i) = [m_1, \Sigma_1, \dots, m_k, \Sigma_k]$$

$$= U^T(y_i, m, \Sigma)$$

$$m = (m_1, m_2, \dots, m_k)$$

$$\Sigma = (\Sigma_1, \dots, \Sigma_k)$$

$$\Rightarrow \sum U^T(y_i, m, \Sigma) z_i + V^T(\pi) z z_i.$$

$$\phi(\theta | \theta^{(P)})$$

$$= E[\log P(x|\theta) | y, \theta^{(P)}]$$

$$= E[\sum U^T(y_i, m, \Sigma) z_i + V(\pi) z_i | y, \theta^{(P)}]$$

$$E[z_i | y, \theta^{(P)}]$$

$$z_i = \begin{bmatrix} z_{i1} \\ z_{i2} \\ \vdots \\ z_{ik} \end{bmatrix} \Rightarrow E[z_{ik} | y, \theta^{(P)}] = P[z_{ik} = 1 | y, \theta^{(P)}]$$

$$= P[y_i \in \text{cluster } k | y_i, \theta^{(P)}]$$

x	0	1
P	P	(1-P)

$E[X] = 0 \cdot P + 1 \cdot (1-P) = P[X=1]$

$$= P[z_i = e_k | y_i, \theta^{(P)}]$$

$$E[z_{ik} | y, \theta^{(P)}] = P[z_i = e_k | y_i, \theta^{(P)}]$$

$$\frac{P[y_i | z_i = e_k, \theta^{(P)}] \cdot P[z_i = e_k | \theta^{(P)}]}{\sum_{k=1}^K P[y_i | z_i = e_k | \theta^{(P)}] \cdot P(z_i = e_k | \theta^{(P)})}$$

$$= \frac{P_k(g_i | m_k^{(p)}, \bar{z}_{ik}^{(p)}) \cdot \pi_{ik}^{(p)}}{\sum_{k'=1}^K P_k'(g_i | m_{k'}^{(p)}, \bar{z}_{ik}^{(p)})}$$

$$= E[\bar{z}_{ik} | y, \theta^{(p)}] = \langle \bar{z}_{ik} \rangle = \bar{z}_{ik}^{(p)} = P[g_i \in k^{\text{th}} \text{ cluster} | \text{current parameters est.}]$$

$$\bar{\Phi}(\theta | \theta^{(p)}) = \sum V^T(g_i, m \bar{z}) \bar{z}^{(p)} + V(\pi) \sum \bar{z}_i^{(p)}$$

$$= \sum V_i^T \bar{z}_i^{(p)} + V^T \bar{z}^{(p)}$$

$$M\text{-Step}, \quad \theta^{(p+1)} = \arg \max_{\theta} \bar{\Phi}(\theta | \theta^{(p)})$$

$$\nabla_{\theta} \bar{\Phi}(\theta | \theta^{(p)}) = 0$$

$$\underbrace{\sum (\nabla_{\theta} V_i^T) \bar{z}_i^{(p)} + (\nabla_{\theta} V^T(\pi)) \varepsilon \bar{z}^{(p)}}_{\pi_k > 0 \quad \sum \pi_k = 1} = 0$$

$$\pi_k > 0 \quad \sum \pi_k = 1.$$

$$\left\{ \begin{array}{l} m_k^{(p+1)} = \frac{\sum \bar{z}_{ik}^{(p)} \cdot g_i}{\sum \bar{z}_{ik}^{(p)}} \\ \sum_k^{(p+1)} = \frac{\sum_k (\bar{z}_{ik}^{(p)} - m_k^{(p+1)}) (g_i - m_k^{(p+1)})^T}{\sum \bar{z}_{ik}^{(p)}} \\ \pi_{ik}^{(p+1)} = \frac{\sum \bar{z}_{ik}^{(p)}}{N} \end{array} \right. \begin{array}{l} \xrightarrow{\text{Probability}} \\ \xrightarrow{\# \text{ of points in cluster } k.} \end{array}$$

7) Determine the # of clusters

for k-means and ~~EM based model~~

① Information Theoretic Criteria.

$$\text{AIC} \quad \text{BIC} \quad \xrightarrow{\text{Bayesian}} \log P(y|\theta) + \text{Penalty}(\theta)$$

\downarrow
XIC

Penalty (θ) = complexity of θ
complexity of model.

② NPB - Non-parametric Bayesian.

$$\text{Parametric } P(y) = P(y|\theta) \xleftarrow{\text{Parameter}}$$

$$\hat{\theta} = \arg \max_{\theta} \log P(y|\theta)$$

$$\text{Non-parametric } P(y) = P(y)$$

$$\hat{P} = \arg \max_P \log P(y)$$

$$\hat{P} = \arg \max_P \left\{ \log P(y) + \log \text{prior}(P) \right\}$$

③ MAP EM

$$\hat{\theta} = \arg \max_{\theta} \left\{ \log P(y|\theta) + \log P(\theta) \right\}$$

Week12

2018年4月16日 17:37

Clustering

Review = EM algorithm

y = observation data \mathcal{D} = D

$$\text{MLF}, \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(y|\theta) \leftarrow \pi P(y|\theta) \\ = \underset{\theta}{\operatorname{argmax}} \log P(D|\theta)$$
$$= \pi \sum_k \pi_k P_k(y|\theta_k)$$

EM algorithm.

x (e.g. $x \in \{y, z\}$) latent hidden variables

E-Step

$$\bar{\phi} = (\theta | \theta^{(p)}) = E[\log P(x|\theta) | y, \theta^{(p)}]$$

(calculating, $E[z | y, \theta^{(p)}]$)

M-step

$$\hat{\theta}^{(p+1)} = \underset{\theta}{\operatorname{argmax}} \bar{\phi}(\theta | \theta^{(p)})$$

Determine K (# of clusters)

MAP EM

MAP EM.

Prior on θ $P(\theta)$

E-step. Same as EM.

$$\text{M-step: } \theta^{(P+1)} = \underset{\theta}{\operatorname{argmax}} \left\{ \bar{\Phi}(\theta | \theta^{(P)}) + \log P(\theta) \right\}$$

Possible prior for θ

final clusters equations for MAP EM

3 clusters

3 labels

Average purity

⊕ priors for mixture model parameters

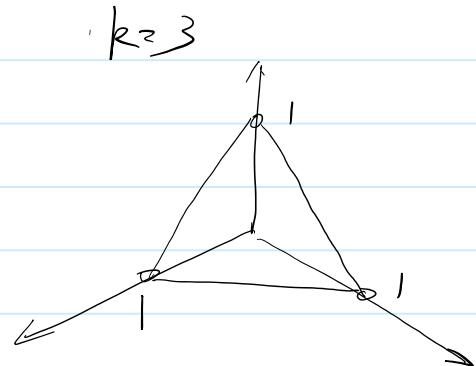
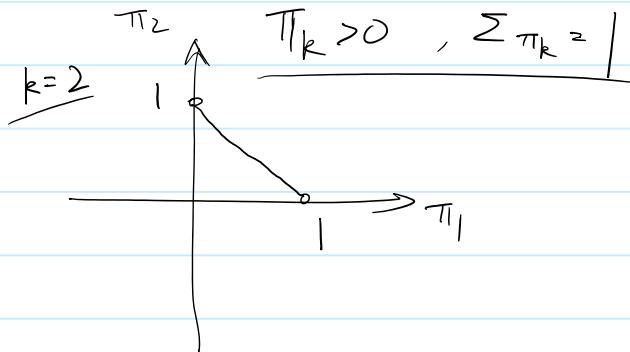
Mixture model.

$$P(y) = \sum_{k=1}^K \pi_k P_k(y | m_k, \Sigma_k)$$

Parameters $\theta = \{\pi_k, [m_k, \Sigma_k]\}$

π_k ← joint prior for $[m_k, \Sigma_k]$

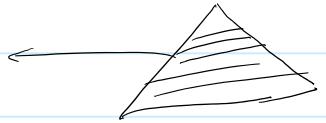
↑
Prior for $\{\pi_k\}$

(5) Prior for $\{\pi_k\}$ 

$$\alpha_1 = \alpha_2 = \alpha_3 = 1$$

Dirichlet Distribution

uniform
density



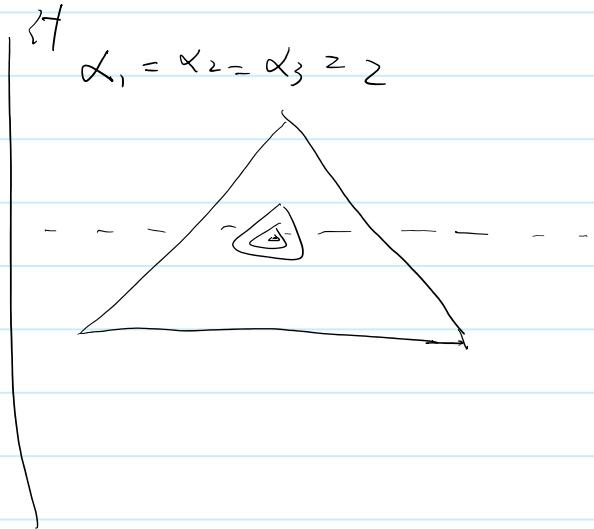
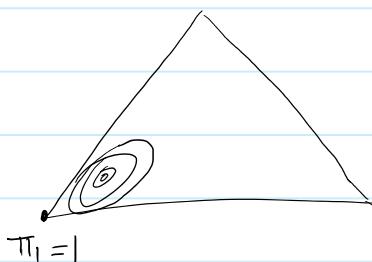
$$\Pi = \{\pi_1, \dots, \pi_k\} \quad \alpha = (\alpha_1, \dots, \alpha_k) > 0$$

$$P(\Pi | \alpha) = \frac{1}{B(\alpha)} \frac{\prod_{k=1}^K \pi_k^{\alpha_k - 1}}{\prod_{k=1}^K \Gamma(\alpha_k)}$$

If $\alpha_k = 1$, all $k \Rightarrow$ uniform distribution

Ref.
B: shop

$$\text{If } k=3, \alpha_1=20, \alpha_2=2, \alpha_3=2$$



⑤ Prior for m_k, Σ_k

$\rightarrow NIN$ distribution

Normal Inverse Wishart

Simpler case, prior for m_k .

Prior for m .

likelihood
Data model

Gaussian $P(y|m, \Sigma)$
mean
cov. matrix

| Prior on m |

$m \sim N(m_0, \Sigma_0)$

constant
hyper
parameters

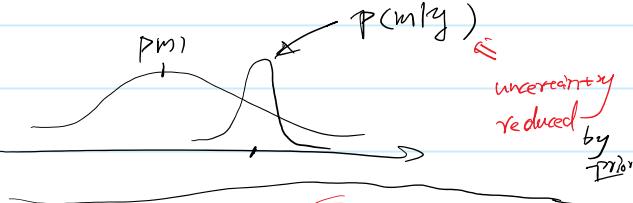
$$\frac{P(\text{Posterior})}{P(m|y)} = \{y_1, y_2, \dots, y_N\}$$

Bayesian inference

$$P(m) \longrightarrow P(m|y)$$

Prior

$$P(y_i|m, \Sigma)$$



$$P(m|y) = \text{Gaussian}$$

$$P(m|y) = \text{Gaussian}$$

Special case $y > 1 \text{ dim.}$

$$P(y|m, \Sigma) = P(y|m, \sigma^2)$$

$$P(m) = P(m|m_0, \Sigma_0)$$

$$P(m|y) = P(m|y, \dots, y_N)$$

$$= P(m|m_N, \sigma_N^2)$$

Special case $y \geq 1$ dim.

$$P(y|m, z) = P(y|m, \sigma^2)$$

$$P(m) = P(m|m_0, \sigma_0^2)$$

$$\begin{aligned} P(m|y) &= P(m|y_1, \dots, y_N) \\ &= P(m|m_N, \sigma_N^2) \end{aligned}$$

$$m_N = \frac{\frac{1}{\sigma^2}}{\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}} \underbrace{\left(\frac{\varepsilon y_i}{N} \right)}_{\hat{m}_m} + \frac{\frac{1}{\sigma_0^2}}{\frac{N}{\sigma^2} + \frac{1}{\sigma_0^2}} m_0$$

$$\frac{1}{\sigma_N^2} = \frac{1}{\sigma_0^2} + \frac{N}{\sigma^2} \rightarrow \sigma_N^2 = \frac{1}{\frac{1}{\sigma_0^2} + \frac{N}{\sigma^2}}$$

When $N \rightarrow +\infty$

$$m_N \rightarrow \hat{m}_m$$

$$\frac{1}{\sigma_N^2} \rightarrow \frac{\sigma^2}{N} 0$$

n -dim case

$$P(y|m, \Sigma)$$

$$P(m|m_0, \Sigma_0)$$

$$P(m|y) = P(m|m_N, \Sigma_N)$$

$$\Sigma_N^{-1} = \Sigma_0^{-1} + N \Sigma^{-1}$$

$$m_N = \Sigma_N^{-1} (\Sigma_N^{-1} \hat{m}_m + \Sigma_0^{-1} m_0)$$

Prior for Σ

$P(y|m, \Sigma)$, Gaussian

↑
fixed.

Wishart inverse distribution.

$P(\Sigma|S_0, V_0) : IW(\Sigma|S_0, V_0)$

↑
prior matrix
for Σ

$$= C |\Sigma|^{-(V_0 + D + 1)/2} \cdot e^{-\frac{1}{2} \text{tr}(S_0 \Sigma^{-1})}$$

↑ trace

$$D = \dim(y)$$

$$= m.$$

$$E[\Sigma] = \frac{S_0}{V_0 - D - 1}$$

$$(V_0 > D + 1)$$

$$A = [a_{ij}]$$

(square)

$$\text{tr}[A] = \sum a_{ii}$$

V_0 - Degree of freedom

$$\text{Property. } \text{trace}(ABC) = \text{trace}(CAB)$$

$$\Rightarrow \text{trace}(X^T \Sigma^{-1} X) \leftarrow H[j]$$

$$= \text{trace}(XX^T \Sigma^{-1})$$

→ A number

↓ ← ↓ []

Posteriors for Σ .

$P(\Sigma|m, y) = IW(\Sigma|S_n, V_n)$

$P(y|m, \Sigma)$, Gaussian.

$$S_n = S_0 + S_m \quad S_m \stackrel{\text{sum}}{\rightarrow} \sum (y_i - m)(y_i - m)^T$$

$$V_n = V_0 + N$$

Prior mean

$$E[\varepsilon] = \frac{s_0}{v_0 - D - 1}$$

Posterior mean

$$E[\varepsilon | m, y] = \int I\!N(\varepsilon | s_n, v_n) \cdot \Sigma d\varepsilon$$

$$= \frac{s_n}{v_n - D - 1} = \frac{s_0 + s_m}{v_0 + N - D - 1} \quad \begin{matrix} \Sigma(y_i - m)(y_i - m)^T \\ \text{circled} \end{matrix}$$

$$= \frac{1}{v_0 + N - D - 1} s_0 + \frac{N}{v_0 + N - D - 1} \cdot \frac{s_m}{N} \quad \begin{matrix} \text{circled} \\ \text{MLE of } \Sigma \end{matrix}$$

When $N \rightarrow \infty$ \Rightarrow MLE of Σ Prior $P(m, \Sigma)$

$$\text{if } P(m, \Sigma) = P(m) \cdot P(\Sigma)$$

 \rightarrow not conjugate wrt $P(y | m, \Sigma)$

Conjugate prior

$$P(m, \Sigma) = P(m | \Sigma) \cdot P(\Sigma)$$

$$= P(m | m_0, \frac{1}{k_0} \Sigma) \cdot I\!N(\Sigma | S_0, V_0) = NIW(m, \Sigma | m_0, k_0, S_0, V_0)$$

Posterior $\xrightarrow{\text{conditional mean.}}$
mode

(MAP estimation)

 \rightarrow Conjugate to Gaussian likelihood.

$$\Rightarrow \text{Posterior } P(m, \Sigma | y) = NIW(m, \Sigma | m_N, k_N, S_N, V_N)$$

$$m_N = \frac{k_0}{k_0 + N} m_0 + \frac{N}{k_0 + N} \left(\frac{\sum y_i}{N} \right)$$

$$k_N = k_0 + N$$

$$V_N = V_0 + N$$

$$S_N = S_0 + S + k_0 m_0 m_0^T + k_N m_N m_N^T$$

MAP estimate.

$$= \underset{m, \Sigma}{\operatorname{argmax}} P(m, \Sigma | y)$$

$$= (m_N, \frac{s_N}{N_N + D + 2})$$

① Applications to clustering
(i.e. mixture model)

\Rightarrow MAP EM.

Recall, $D = \{y_1, y_2, \dots, y_N\} = y$

$$P(y) = \sum \pi_k P(y | m_k, \Sigma_k)$$

$$= P(y|\theta)$$

$$\theta = \{\pi_k, m_k, \Sigma_k\}$$

Complete data. $X = (y, z)$

z = hidden variables = $\{z_1, z_2, \dots, z_N\}$

$z_i = e_k$ if $y_i \in k^{\text{th}}$ cluster.

\uparrow
| host vector

$$= \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix}$$

$$\text{E-step} \rightarrow \bar{\phi}(\theta|\theta^{(P)})$$

$$= E[\log \hat{P}(x|\theta) | y, \theta^{(P)}]$$

$$= \sum_k \left(\bar{z}_{ik} \log \pi_k + \bar{z}_{ik} \log P(y_i | m_k, \Sigma_k) \right)$$

\uparrow

$E[z_{ik} | y, \theta^{(P)}]$

M-step (MAP, M-step)

$$\theta^{(P+1)} = \arg \max \left\{ \bar{\phi}(\theta|\theta^{(P)}) + \log \hat{P}(\theta) \right\}$$

$$\log \hat{P}(\theta) = \log P(\pi) + \sum \log P(m_k, \Sigma_k)$$

$$\left\{ \pi_1, \pi_2, \dots, \pi_k \right\}$$

Priechet

$$NIW(m_k, \Sigma_k | m_o, k_o, V_o, S_o)$$

$$m_k^{(P+1)} = \frac{\bar{z}_k}{\bar{z}_k + k_o} m_{k,ML} + \frac{k_o}{\bar{z}_k + k_o} m_o$$

$$\bar{z}_k = \sum_{k=1}^N \bar{z}_{ik} \frac{\sum_i z_{ik} y_i}{\sum_i z_{ik}}$$

$$\Sigma_k = \frac{s_o + s_k}{n_o + \bar{z}_k + D + 2} + \frac{k_o \bar{z}_k}{k_o + \bar{z}_k} (m_{k,ML} - m_o) (m_{k,ML} - m_o)^T$$

$$s_k = \sum_i (y_i - m_{k,ML}) (y_i - m_{k,ML})^T$$

↑
Sum

$$\pi_k = \frac{\bar{z}_k + \alpha_{k-1}}{N + \sum_k \alpha_k - K}$$

Week13

2018年4月16日 18:46

7. Variational Influence.

1) Why?

influence:

$$\rightarrow \text{MLE}, \hat{\theta} = \arg \max_{\theta} \log P(D|\theta)$$

$$\begin{array}{c} \text{observation} \\ \downarrow \\ \log P(y|\theta) \\ \downarrow \\ \log P(D|\theta) \end{array}$$

$$\rightarrow \text{MAP} \quad \hat{\theta} = \arg \max_{\theta} \left\{ \log P(D|\theta) + \log P(\theta) \right\}$$

$$\rightarrow \text{posterior} \quad P(z|y, \theta)$$

↑
hidden variable ↓
observed variable.

problem.

analytic derivation can be difficult.

→ influence. difficult.

Solution.

EM algorithm (require $P(z|y, \theta)$)

$$P(z|y, \theta^*)$$

can be calculated for mixture case.

However, many cases can't be calculated.

Solution = Variational influence.

Basic idea of Variational/ Influence.

P difficult to work with.

\Rightarrow find g , good appr. to P .

work with g .

① How to measure the difference between P & g .

② How to find the closest g for P ?

③ Measure of difference of two $\frac{\text{pdfs}}{\text{continuous}}$ / $\frac{\text{prob. distribution}}{\text{discrete}}$

$$P = P(x) \text{ , pdf.}$$

$$g = g(x) \text{ , pdf.}$$

$$d(P, g) = \text{SQE} = C \int_{-\infty}^{+\infty} (P - g)^2(x) dx$$

$$\text{also, } C \int_{-\infty}^{+\infty} (P - g)^2 p(x) dx$$

$p(x)$

ave code length
for data comes from P but coded as if it comes g $H(P, g) - H(g)$ for data with P

Directed

KL divergence

$$D(P//g) = \int_{-\infty}^{+\infty} P(x) \log \frac{P(x)}{g(x)} dx = \int_{-\infty}^{+\infty} P(x) (\log P(x) - \log g(x)) dx \geq 0$$

$$D(g//P) = \int_{-\infty}^{+\infty} g(x) \log \frac{g(x)}{P(x)} dx$$

3) Variational (Calculus of variation)

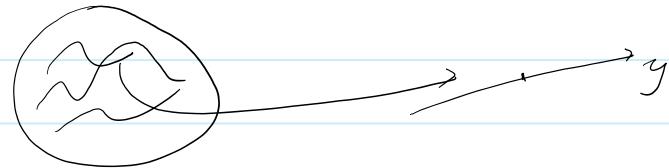
function, $y = f(x)$ (e.g. $y = f(x) = x^2$)

mapping $\mathbb{R} \rightarrow \mathbb{R}$

Δ functional, $F(g) = y$

function \nearrow y \uparrow
 g number

mapping: Function space $\longrightarrow \mathbb{R}$



$$F(g) = \int_{-\infty}^{+\infty} |g(x)|^2 dx.$$

\uparrow
 $g = g(x)$

operator,

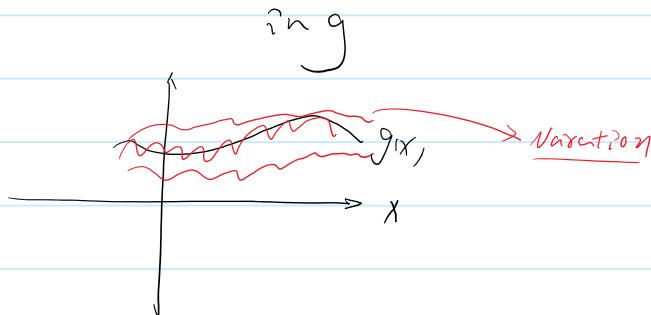
mapping function \longrightarrow function.

$$H(g) = (g * h)(x)$$

Optimize a function $\frac{df}{dx} = 0$ (e.g. $f(x) = x^2$, $\frac{df}{dx} = 2x = 0 \Rightarrow x=0$)

Optimize a functional $\int g$

Small variation in F Small variation in g



Special case

$$\delta f(x) = g_\theta(x) \quad (\text{e.g. } g_{\text{fit}} = g_\theta(x))$$

$$\frac{\delta g}{\delta \theta} = \nabla_\theta F = N(x / m \cdot \sigma^2)$$

Example for special case,

X : random variable

Y : random variable.

$\hat{Y} = g(x) = \text{function estimation of } Y$

$$= g_{\text{fit}}(x) = ax + b$$

$$g^* = \underset{g}{\operatorname{arg\,min}} E[(Y - g(x))^2] \Rightarrow \text{min a functional.}$$

add constraints

$$= \underset{a, b}{\operatorname{arg\,min}} E[\underbrace{(Y - (ax+b))^2}_{\varepsilon}] \quad \frac{\partial \varepsilon}{\partial a} = 0, \frac{\partial \varepsilon}{\partial b} = 0$$

A more general case.

$$F(g) = \int L(x, g(x), g'(x)) dx$$

$$\frac{\delta F}{\delta g} = 0 \quad \frac{\partial L}{\partial g} - \frac{d}{dx} \frac{\partial L}{\partial g'} = 0$$

$$\Rightarrow g(x)$$

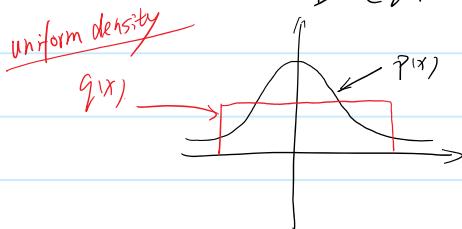
Example, x , random variable.

$$p(x) \sim N(0, \sigma^2) \quad \begin{matrix} \rightarrow \text{fixed} \\ \text{variance} \end{matrix}$$

$$g(x) \sim U[-a, +a], a > 0$$

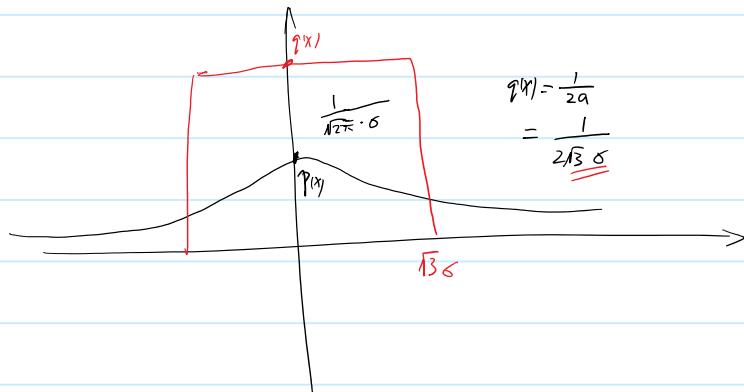
find $g(x)$ such that

$D(G \nparallel P)$ is minimized.



$$\begin{aligned} D(P \parallel Q) &= \int_{-\infty}^{+\infty} q(x) \log \frac{q(x)}{p(x)} dx \\ &= \int_{-\infty}^{+\infty} \left[\frac{1}{2\sigma} e^{-\frac{x^2}{2\sigma^2}} \right] \log \frac{\frac{1}{2\sigma} e^{-\frac{x^2}{2\sigma^2}}}{\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}} dx \\ &= C - \log a + \frac{a^2}{6\sigma^2} \end{aligned}$$

$$\begin{aligned} \frac{d}{da} &= 0 \Rightarrow -\frac{1}{a} + \frac{a}{3\sigma^2} = 0 \\ \Rightarrow a^2 &= 3\sigma^2 \\ a &= \sqrt{3}\sigma \end{aligned}$$



In practice

$P(\cdot)$: real distribution

$Q(\cdot)$: approximation, easier to work with

$$\Delta Q = \underset{Q}{\operatorname{argmin}} D(P \parallel Q), \quad Q = \underset{Q}{\operatorname{argmin}} D(Q \parallel P) \Rightarrow E \left[Q \cdot \log \frac{Q}{P} \right]$$

Better approximation
hard to optimize

easier to optimize

4). Variational lower bound

Application: interested in a posterior:

x — data, observable.

z — hidden/latent variable.

$$ML_F = \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log P(x|\theta)$$

hard to work with.

$$\text{Posterior}, P(z|x, \theta) \quad \leftarrow \text{hard to derive or}$$

Week14

2018年4月30日

18:16

2

2018年4月30日 13:37

3

2018年4月30日 13:38

4

2018年4月30日 13:37

5

2018年4月30日 13:38

6

2018年4月30日 13:37

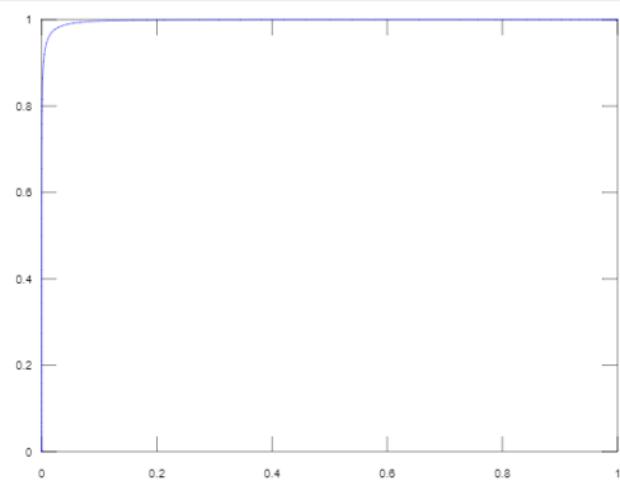
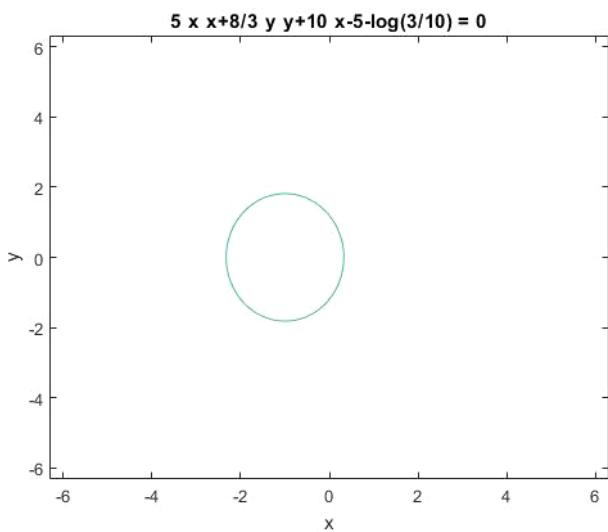
7

2018年4月30日 13:38

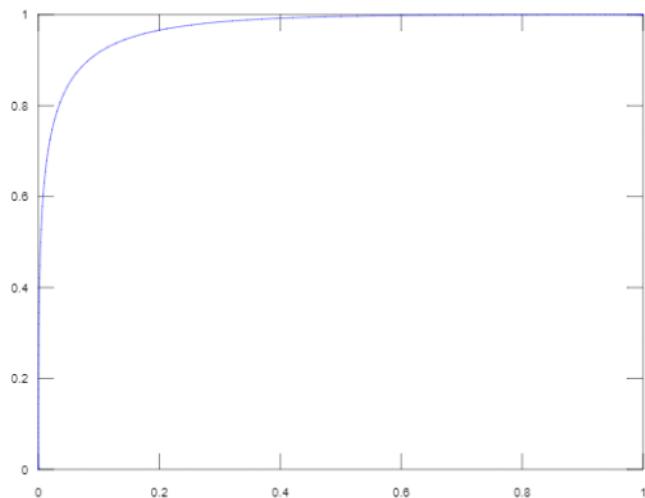
8

2018年4月30日 13:38

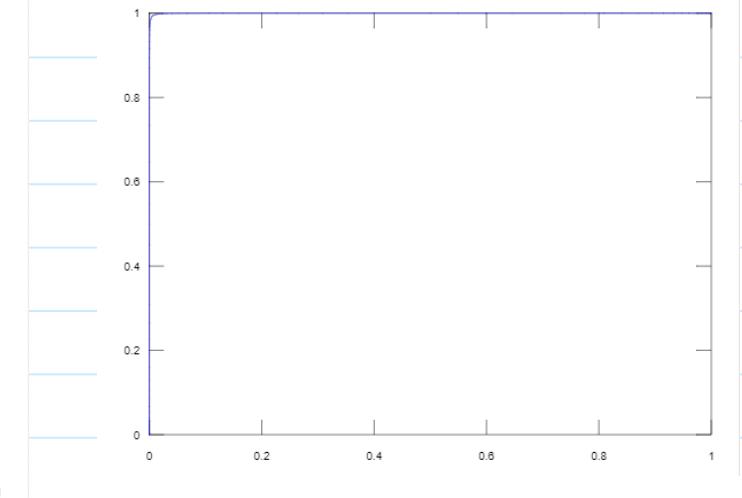
Appendix



Sigma=1



Sigma=1.5



Sigma=0.8